

Soft Random Graphs in Probabilistic Metric Spaces & Inter-graph Distance

Kangrui Wang* and Dalia Chakrabarty^{†,‡},

[‡] *Corresponding author*
Department of Mathematics
Brunel University London
Uxbridge
Middlesex UB8 3PH
U.K.

dalia.chakrabarty@brunel.ac.uk

Abstract: We present a new method for learning Soft Random Geometric Graphs (SRGGs), drawn in probabilistic metric spaces, with the connection function of the graph defined as the marginal posterior probability of an edge random variable, given the correlation between the nodes connected by that edge. In fact, this inter-node correlation matrix is itself a random variable in our learning strategy, and we learn this by identifying each node as a random variable, measurements of which comprise a column of a given multivariate dataset. We undertake inference with Metropolis with a 2-block update scheme. The SRGG is shown to be generated by a non-homogeneous Poisson point process, the intensity of which is location-dependent. Given the multivariate dataset, likelihood of the inter-column correlation matrix is attained following achievement of a closed-form marginalisation over all inter-row correlation matrices. Distance between a pair of graphical models learnt given the respective datasets, offers the absolute correlation between the given datasets; such inter-graph distance computation is our ultimate objective, and is achieved using a newly introduced metric that resembles an uncertainty-normalised Hellinger distance between posterior probabilities of the two learnt SRGGs, given the respective datasets. Two sets of empirical illustrations on real data are undertaken, and application to simulated data is included to exemplify the effect of incorporating measurement noise in the learning of a graphical model.

AMS 2000 subject classifications: Graphical methods, 62xx; Random graphs, 05C80; Measures of association (correlation, canonical correlation, etc), 62H20; Distance in graphs, 05C12.

Keywords and phrases: Soft random geometric graphs, Probabilistic metric spaces, Inter-graph distance, Hellinger distance, Metropolis by block update, Human disease-symptom network.

1. Introduction

Graphical models of complex multivariate datasets, manifest intuitive illustrations of the correlation structures of the data, and are of interest in different disciplines (Airoldi, 2007; Bandyopadhyay and Canale, 2016; Benner et al., 2014;

*Alan Turing Institute

[†]Department of Mathematics, Brunel University London

Carvalho and West, 2007; Whittaker, 2008). Much work is present in the Statistics literature on the intrinsic correlation structure of a multivariate dataset that comprises multiple measurements of a vector-valued observable, where it is common practice to model the joint probability distribution of a set of such observable values, as matrix-Normal (Gruber and West, 2016; Ni et al., 2017; Wang and West, 2009).

In this paper, we present a method for learning the correlation structure of a multivariate dataset, and its graphical model. The general method is advanced for iterative Bayesian learning of the correlation matrix, at each update of which, a soft random geometric graph (SRGG) (Giles et al., 2016; Penrose, 2003, 2016) of the data is updated, where any such SRGG is drawn in probabilistic metric space, (Menger, 1942; Schweizer and Sklar, 1983), such that its connection function is the location-dependent marginal posterior probability of an edge, given the correlation between the nodes that straddle this edge, and the chosen cutoff radius is a probability in the probabilistic metric space that the SRGG is drawn in. In fact, such an SRGG is shown to be underlined by a non-homogeneous Poisson process with an intensity that is dependent on the location, or more precisely, the nodes, and thereby on the inter-nodal correlation. Thus, the point process that generates this SRGG is compounded with the process that generates the matrix-valued correlation variable. The full graphical model of the data is then defined using the sequence of SRGGs generated across iterations of this learning scheme. In this method, inference on uncertainties of both the correlation matrix and graphical model is undertaken, and we can acknowledge measurement errors in learning both random structures as well.

The learning of the graph and correlation structure are undertaken Bayesianly, using Bayesian inference that is implemented via Markov Chain Monte Carlo (MCMC) inference techniques; to be precise, a Metropolis-with-2-block-update-based algorithm is implemented (Robert and Casella, 2004), to make inference on the correlation matrix given the data, and on the SRGG given the updated correlation. However, for learning large networks, for which such iterative inference is not practical, we modulate the method to accommodate this concern. Then we undertake the network learning as a single SRGG.

The ultimate aim behind our learning of the graphical model of a given data, is to compute the distance between the graphical models learnt for a pair of given datasets, in order to thereby compute the strength of the inter-data correlation. We compute the distance between the graphical models learnt for the respective dataset, using a new metric δ that we introduce (in Theorem 4.1), where this metric is akin to the Hellinger distance between the posterior probabilities of the graphical models given the respective correlation structure of the given datasets, normalised by the uncertainties in each of the learnt graphical models. Such distance informs on the absolute correlation between the pair of multivariate datasets, for which the graphical models are learnt.

Objective and comprehensive uncertainties on the Bayesianly learnt graphical model of given multivariate data, are sparsely available in the literature. Such uncertainties can potentially be very useful in informing on the range of models that describe the partial correlation structure of the data at hand. Madigan and

Raftery (1994) discuss a method for computing model uncertainties by averaging over a set of identified models, and they advance ways for the computation of the posterior model probabilities, by taking advantage of the graphical structure, for two classes of considered models, namely, the recursive causal models (Kiiveri et al., 1984) and the decomposable log-linear models (Goodman, 1970). This method allows them to select the “best models”, while accounting for model uncertainty. Our method on the other hand, provides a direct and well-defined way of learning uncertainties of the graphical model of a given multivariate data.

However, we wish to extend such learning to higher-dimensional data, for example, to a dataset that is cuboidally-shaped, given that it comprises multiple measurements of a matrix-valued observable. Hoff (2011); Xu et al. (2012); Wang & Chakrabarty (<https://arxiv.org/abs/1803.04582>), advance methods to learn the correlation in high-dimensional data in general. For a rectangularly-shaped multivariate dataset, the pioneering work by Wang and West (2009) allows for the learning of both the inter-row and inter-column covariance matrices, and therefore, of two graphical models. Ni et al. (2017) extend this approach to high-dimensional data. However, a high-dimensional graph showing the correlation structure amongst the multiple components of a general hypercuboidally-shaped dataset, is not easy to visualise or interpret. Instead, in this paper, we treat the high-dimensional data as built of correlated rectangularly-shaped data slices, given each of which, the inter-column (partial) correlation structure and graphical model are Bayesianly learnt, along with uncertainties, subsequent to our closed-form marginalisation over all inter-row correlation matrices (in Section 2.5, unlike in the work of Wang and West (2009)). By invoking the uncertainties learnt in the graphical models, we advance a new inter-graph distance metric (Section 4), based on the Hellinger distance (Banerjee et al., 2015; Matusita, 1953) between the posterior probability densities of the pair of graphical models that are learnt given the respective pair of such rectangularly-shaped data slices. We use a corresponding affinity measure to then infer on the correlation between the pair of datasets (Section 4.1), permitting the correlation structure of the high-dimensional dataset thereby. Thus, by computing the pairwise inter-graph distance between each learnt pair of graphs, we can avoid the inadequacy of trying to capture spatial correlations amongst sets of multivariate observations, by “computing partial correlation coefficients and by specifying and fitting more complex graphical models”, as was noted by Guinness et al. (2014). An additional advantage is that our method offers the inter-graph distance for two differently sized datasets.

That we do not learn the graphical model as an Erdos-Renyi graph is because we wish to utilise the fully Bayesian nature of the inference that we implement here, without resorting to any unrealistic assumptions – such as decomposability. Effectively, we wish to learn the uncertainty-included graphical model of noisy data, as distinguished from making inference on the graph (i.e. writing its posterior) clique-by-clique. Also, we do not need to be reliant on the closed-form nature of the posteriors to sample from, i.e. we do not need to invoke conjugacy to affect our learning. Indeed, to contextualise to a common practice in Bayesian learning of undirected graphs, a Hyper-Inverse-Wishart prior is typ-

ically imposed on the covariance matrix of the data, as this then allows for a Hyper-Inverse-Wishart posterior of the covariance, which in turn implies that the marginal posterior of any clique is Inverse-Wishart – a known, closed-form density (Dawid and Lauritzen, 1993; Lauritzen, 1996). Inference is then rendered easier, than when generating posterior samples from a non-closed form posterior, (using techniques such as MCMC). Now, if the graph is not decomposable, and a Hyper-Inverse-Wishart prior is placed on the covariance matrix, the resulting Hyper-Inverse-Wishart joint posterior density that can be factorised into a set of Inverse-Wishart densities, cannot be identified as the clique marginals. Expressed differently, the clique marginals are not closed-form when the graph is not decomposable. However, this is not a worry in our learning as we can undertake our learning irrespective of the validity of decomposability. Our pursued graphical model is a soft RGG that is drawn in a (normed) probabilistic space, where the location-dependent affinity measure between a pair of vertices of the graph is conditional on the correlation learnt between the pair of random variables at these two vertices.

This paper is organised as follows. The following section deliberates upon the methodology development that we advance towards the learning of the SRGG and the inter-column correlation matrix of a given dataset. In the subsequent section, issues relevant to the inference on the unknowns is discussed, along with definition of uncertainties in the learnt graphical model. The metric used for computing the inter-graph distance, is then discussed in Section 4, while Section 5 presents our modulated learning methodology to accommodate challenges of learning large networks as SRGGs. Section 6 presents the empirical illustration on 2 real datasets and comparison against existing results, while distance computation between the learnt graphical models of these 2 data, is discussed in Section 7. In Section 8, we learn the large disease-phenotype network, and compare our results with those reported earlier (Hoehndorf et al., 2015). The paper is rounded up with a section that summarises the main findings and the conclusions.

The attached Supplementary Material elaborates on certain aspects of our work. This includes quantitative comparison of the results that we obtain using the real datasets that we illustrate our methodology upon (Sections 7 and 9 of the Supplementary Material, along with outputs of such a comparative exercise included in Figures 12, 13, 14 of the Supplement), with results that are available in the literature, or obtained independently by us. Importantly, detailed model checking is discussed in Section 5 of the Supplementary Material, in the context of an empirical illustration made to a simulated dataset (presented in Section 4 of the Supplement). Convergence signatures of our MCMC chains are borne by the extra results that are included in Figures 9, 10, and 11 of the Supplement, in addition to discussions below in Section 6.

2. Background

Let points X_1, \dots, X_p be independent, with the random variable $X_i \in \mathcal{X} \subseteq \mathbb{R}^d$ be s.t. $X_i \sim f(\theta_1^{(i)}, \dots, \theta_q^{(i)})$, with the parameters of the *pdf* of X_i given as $\theta_1^{(i)}, \dots, \theta_q^{(i)} \in \mathbb{R}$.

The d -dimensional, soft random geometric graph (SRGG) $\mathcal{G}_\phi(\mathbf{V})$ on the vertex set $\mathbf{V} := \{X_1, \dots, X_p\}$, with each of the p points $X_1, \dots, X_p \in \mathcal{X} \subseteq \mathbb{R}^d$ assigned a random coordinate in the box $[0,1]^d$, is s.t. probability of edge G_{ij} between the i -th and j -th nodes ($i \neq j; i, j \in \{1, \dots, p\}$), is given by a function $\phi(\cdot)$ of the distance between point X_i and point X_j . Here $\phi : \mathcal{X} \rightarrow [0, 1]$ is referred to as the connection function, following Penrose (2016).

2.1. Probabilistic metric space: background

We draw our SRGG in a probabilistic metric space \mathcal{X} (Menger, 1942), s.t. to any $X_1, X_2 \in \mathcal{X}$, a probability distribution $F_{X_1, X_2}(X)$ is assigned, where $F_{X_1, X_2}(0) = 0$ (similar to the assignment of a non-negative number to any two points in a metric space). Let all distributions that abide by the constraint $F_{\cdot, \cdot}(0) = 0$, live in space $\mathcal{F}_+ \subset \mathcal{F}$, where probability distributions live in \mathcal{F} . We note that for $X_1 = X_2$, the distribution $F_{X_1, X_2}(\cdot) = \ell_0(\cdot)$, where the distribution $\ell_a(\cdot) \in \mathcal{F}_+$ is s.t. $\ell_a(x) = 0$ if $x \leq 0$ and $\ell_a(x) = 1$ if $x > 0$, for $a \in \mathbb{R}_{\geq 0}$.

Definition 2.1. A probabilistic metric space is the triple

$$\{\mathcal{X}, F, \Delta\},$$

where the probabilistic distance is $F_{X_1, X_2} : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{F}_+$, $\forall X_1, X_2 \in \mathcal{X}$, with

$$F_{X_1, X_2}(\cdot) = \ell_0(\cdot) \iff X_1 = X_2;$$

$$F_{X_1, X_2}(\cdot) = F_{X_2, X_1}(\cdot), \forall X_1, X_2 \in \mathcal{X}; \text{ and}$$

$$\text{triangle function } \Delta \text{ is s.t. } F_{X_1, X_3}(\cdot) \geq \Delta(F_{X_1, X_2}, F_{X_2, X_3})(\cdot), \forall X_1, X_2, X_3 \in \mathcal{X},$$

with the triangle function Δ defined as a binary operation on \mathcal{F}_+ , with respect to the triangle norm D s.t.

$$\Delta_D(F_{X_1, X_2}, F_{X_3, X_4})(x) = \sup[D(F_{X_1, X_2}(u), F_{X_3, X_4}(v)); u + v = x],$$

where the triangle norm D (Fodor, 2004) is defined as the binary operation on

$$\text{interval } [0, 1], \text{ s.t. } D(y_1, y_2) = D(y_2, y_1) \forall y_1, y_2 \in [0, 1]; D(y_1, y_2) \geq D(y_3, y_4) \forall y_1 \geq$$

$$y_3; y_2 \geq y_4; y_1, y_2, y_3, y_4 \in [0, 1]; D(y_1, D(y_2, y_3)) = D(D(y_1, y_2), y_3)$$

$$\forall y_1, y_2, y_3 \in [0, 1]; \text{ and } D(y, 1) = y, \forall y \in [0, 1],$$

2.2. Marginal posterior of edge parameter as connection function

We draw our SRGG in a probabilistic metric space X .

Definition 2.2. Points X_1, \dots, X_p are assigned random coordinates in a 2-dimensional square, to construct a 2-dimensional SRGG on the vertex set $\mathbf{V} = \{X_1, \dots, X_p\}$, with the probability of the edge between X_i and X_j (where $X_i \neq X_j; X_i, X_j \in \mathbf{V}$) given by the marginal posterior probability of the random edge variable G_{ij} , conditional on the (partial) correlation ρ_{ij} between the random variables X_i and X_j , i.e. the connection function in our SRGG is the marginal posterior of G_{ij} , given ρ_{ij} :

$$m(G_{ij} | \rho_{ij}).$$

Remark 2.1. In our SRGG, the connection function is the affinity measure between a pair of points, defined by the marginal posterior of the edge variable connecting them in a (normed) probabilistic metric space (Menger, 1942). This affinity measure is complementary to the distance between this pair of points in this space.

An immediate definition of the posterior probability density of our SRGG would be the joint posterior probability of the edge parameters ($\{G_{ij}\}_{i \neq j; i, j=1}^p$) given the partial correlation structure, as:

$$\pi(G_{11}, G_{12}, \dots, G_{p, p-1} | \mathbf{R}) \propto \ell(G_{11}, G_{12}, \dots, G_{p, p-1} | \mathbf{R}) \pi_0(G_{11}, G_{12}, \dots, G_{p, p-1}),$$

where $\pi_0(G_{11}, G_{12}, \dots, G_{p, p-1})$ is the prior probability density on the edge parameters $\{G_{ij}\}_{i \neq j; i, j=1}^p$, and $\ell(G_{12}, \dots, G_{1p}, G_{23}, \dots, G_{2p}, G_{34}, \dots, G_{p, p-1} | \mathbf{R})$ is the likelihood of the edge parameters, given the partial correlation matrix $\mathbf{R} = [\rho_{ij}]$.

We choose a prior on G_{ij} that is *Bernoulli*(0.5) $\forall i, j$, i.e.

$$\pi_0(G_{11}, G_{12}, \dots, G_{p, p-1}) = \prod_{i, j=1; i \neq j}^p 0.5^{G_{ij}} 0.5^{1-G_{ij}};$$

thus, the prior is independent of the edge parameters. In applications marked by more information, we can resort to stronger priors. However, as we will soon see, this posterior definition needs updating.

We choose to define this likelihood as a function of the (squared) Euclidean distance between the ‘‘observation’’, (i.e. the absolute value of ρ_{ij}), and the unknown parameter G_{ij} , with the squared distance normalised by a squared scale length, i.e. the variance parameter v_{ij} , for all relevant ij -pairs. Thus, the

unknown parameters in the model are the edge parameters $\{G_{ij}\}_{i \neq j; i, j=1}^p$ and variance parameters $\{v_{ij}\}_{i \neq j; i, j=1}^p$. In light of these newly introduced variance parameters, we rewrite our likelihood (of the unknown edge and variance parameters, given the partial correlation matrix), as:

$$\ell(G_{12}, \dots, G_{1p}, G_{23}, \dots, G_{2p}, \dots, G_{p,p-1}, v_{12}, \dots, v_{1p}, v_{23}, \dots, v_{2p}, \dots, v_{p,p-1} | \mathbf{R}).$$

Now, we choose a parametric form of the likelihood, given its anticipated properties. Likelihood increases (decreases) as distance between $|\rho_{ij}|$ and G_{ij} decreases (increases). Also the likelihood is invariant to change of sign of $|\rho_{ij}| - G_{ij}$. Given this, we model our likelihood of the edge and variance parameters, given \mathbf{R} as

$$\begin{aligned} \ell(G_{12}, \dots, G_{1p}, G_{23}, \dots, G_{2p}, \dots, G_{p,p-1}, v_{12}, \dots, v_{1p}, v_{23}, \dots, v_{2p}, \dots, v_{p,p-1} | \mathbf{R}) \\ = \prod_{i \neq j; i, j=1}^p \frac{1}{\sqrt{2\pi v_{ij}}} \exp \left[-\frac{(G_{ij} - |\rho_{ij}|)^2}{2v_{ij}} \right], \end{aligned} \quad (2.1)$$

where the variance parameters $\{v_{ij}\}_{i \neq j; i, j=1}^p$ are indeed hyper-parameters that are also learnt from the data. The variance parameters are assigned uniform prior probability in the interval $[0, 1]$.

Thus, the joint posterior probability of the edge and variance parameters is

$$\pi(\{G_{ij}\}_{i \neq j; i, j=1}^p, \{v_{ij}\}_{i \neq j; i, j=1}^p | \mathbf{R}) = K \prod_{i \neq j; i, j=1}^p \frac{1}{\sqrt{2\pi v_{ij}}} \exp \left[-\frac{(G_{ij} - |\rho_{ij}|)^2}{2v_{ij}} \right], \quad (2.2)$$

where $K > 0$ is a constant that incorporates information on the ij -independent Bernoulli(0.5) prior, Uniform[0,1] prior, and probability of \mathbf{R} that defines the denominator of Bayes rule.

Then the marginal posterior probability of the ij -th edge parameter G_{ij} , given the partial correlation ρ_{ij} between X_i and X_j , is:

$$\begin{aligned} m(G_{ij} | \rho_{ij}) &= K \int_0^1 \frac{1}{\sqrt{2\pi v}} \exp \left[-\frac{(G_{ij} - |\rho_{ij}|)^2}{2v} \right] dv \\ &= K \left[\sqrt{\frac{2}{\pi}} \exp \left(-\frac{(G_{ij} - |\rho_{ij}|)^2}{2} \right) - (|G_{ij} - |\rho_{ij}||) \operatorname{erfc} \left(\frac{|G_{ij} - |\rho_{ij}||}{\sqrt{2}} \right) \right], \end{aligned} \quad (2.3)$$

where $\operatorname{erfc}(\cdot)$ is the complementary error function.

Theorem 2.1. *Setting the global constant $K = 1$, separation D_{ij} between points $X_i \in \mathcal{X}$ and $X_j \in \mathcal{X}$ is a norm in the probabilistic metric space \mathcal{X} , where*

$$\begin{aligned} D_{ij} &:= m(G_{ij} | \rho_{ij}) + |G_{ij} - |\rho_{ij}||, \\ &= \left[\sqrt{\frac{2}{\pi}} \exp \left(-\frac{(G_{ij} - |\rho_{ij}|)^2}{2} \right) + (|G_{ij} - |\rho_{ij}||) \operatorname{erf} \left(\frac{|G_{ij} - |\rho_{ij}||}{\sqrt{2}} \right) \right], \end{aligned} \quad (2.4)$$

where we have recalled that the error function $\text{erf}(\cdot) = 1 - \text{erfc}(\cdot)$, and $m(G_{ij}|\rho_{ij})$ is defined in Equation 2.3.

The proof of this theorem is in Section 1 of the attached Supplementary Materials.

The global scale K is chosen to ensure that D_{ij} is a metric in the probabilistic metric space \mathcal{X} .

Figure 1 shows the variation with $|G_{ij} - |\rho_{ij}||$, in the unscaled distance D_{ij} between the i -th and j -th nodes; the unscaled affinity measure $m(G_{ij}|\rho_{ij})$ between points X_i and X_j ; and the difference between the unscaled distance and the affinity measure, for G_{ij} set to 1 (trends displayed in the left panel), and G_{ij} set to 0 (results shown on the right). In our SRGG, the ij -th edge exists with the probability that is given by the probability that affinity $m(G_{ij}|\rho_{ij})$ between X_i and X_j exceeds cutoff probability τ . Now, probability for the ij -th edge to exist, increases with increase in $|\rho_{ij}|$; so we expect $m(G_{ij} = 1|\rho_{ij})$ to increase with $|\rho_{ij}|$. This trend is borne in the results displayed in Figure 1. From this figure, we also see that the affinity measure defined in terms of the marginal posterior of the edge, given the partial correlation, is complementary to the distance D_{ij} , in the sense that as affinity increases, distance decreases, for a given value of G_{ij} , $\forall i, j \in \{1, 2, \dots, p\}; i \neq j$.

2.3. Edge set of our SRGG

Definition 2.3. Then edge G_{ij} exists in the graph, independently of any other edge, if and only if, affinity between X_i and X_j exceeds a threshold probability τ , i.e.

$$m(G_{ij}|\rho_{ij}) \geq \tau.$$

Thus, the value of the ij -th edge parameter G_{ij} is

$$\begin{aligned} g_{ij} &= 1 && \text{if } m(G_{ij}|\rho_{ij}) \geq \tau; i \neq j; i, j \in \{1, \dots, p\} \\ g_{ij} &= 0 && \text{otherwise} \end{aligned} \quad (2.5)$$

Here the ij -th edge exists in the graph if $g_{ij} = 1$, and does not exist if $g_{ij} = 0$.

Thus, the edge set of our graph is

$$E_p = \{G_{ij} : m(G_{ij}|\rho_{ij}) \geq \tau; X_i \neq X_j; X_i, X_j \in \mathbf{V}\}.$$

Remark 2.2. As ρ_{ij} is dependent on the i -th and j -th vertices, marginal $m(G_{ij}|\rho_{ij})$ that is a function of ρ_{ij} , is i and j -dependent too, implying that the connection

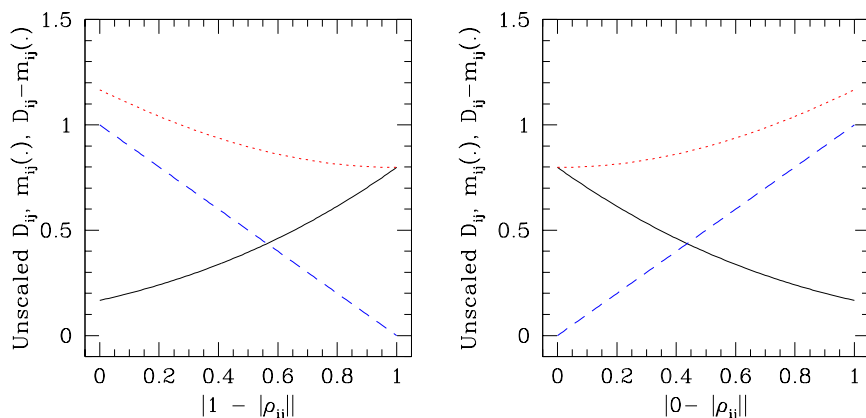


FIG 1. Figure displaying variation with the input $|G_{ij} - |\rho_{ij}||$, of unscaled distance D_{ij} between the i -th and j -th nodes (in dotted lines); the unscaled affinity measure $m(G_{ij}|\rho_{ij})$ between points X_i and X_j (in solid line); and the difference between the unscaled distance and the affinity measure (in dashed lines), for $G_{ij} = 1$ (left panel) and $G_{ij} = 0$ (right panel). Here the input variable is the absolute of the difference between value (0 or 1) of the edge variable G_{ij} between i -th and j -th nodes, and the absolute partial correlation $|\rho_{ij}|$ between points X_i and X_j . As $|\rho_{ij}|$ increases, probability for the edge connecting X_i and X_j to exist – given by the affinity measure $m_{ij}(G_{ij} = 1|\rho_{ij})$ between points X_i and X_j in our SRGG to exceed a chosen cutoff probability – increases, while the probability for this edge to not exist (i.e. for G_{ij} to be 0) decreases, as shown on the right. Distance between points X_i and X_j is defined to be complimentary to the affinity.

function is location-dependent in this graph. Then following Penrose (2016), our notation for this location-dependent SRGG defined at threshold parameter τ is $\mathcal{G}_{m, \mathbf{R}}(\mathbf{V}, \tau)$, where the partial correlation matrix is

$$\mathbf{R}^{(p \times p)} = [\rho_{ij}],$$

where cardinality of \mathbf{V} is p .

2.4. Point process background for our SRGG

In this sub-section, we present the underlining Point Process (PP) that can generate our Soft RGG $\mathcal{G}_{m, \mathbf{R}}(\mathbf{V}, \tau)$. To undertake this, we now model the points $X_i, X_j \in \mathcal{X}$, as random variables X_i and X_j that are Normally distributed with means μ_i and μ_j , and the variance σ^2 .

Theorem 2.2. Distance between random variables $X_i, X_j \in \mathcal{X}$ – where \mathcal{X} is

the host probabilistic metric space – is given as:

$$D(X_i, X_j) = \frac{2\sigma}{\sqrt{\pi}} \exp\left(\frac{-|\mu_i - \mu_j|^2}{4\sigma^2}\right) + |\mu_i - \mu_j| \operatorname{erf}\left(\frac{|\mu_i - \mu_j|}{2\sigma}\right),$$

with $X_i \sim \mathcal{N}(\mu_i, \sigma^2)$, $X_j \sim \mathcal{N}(\mu_j, \sigma^2)$.

Proof. From definition of distance in probabilistic metric spaces, it follows that distance between X_i and X_j in \mathcal{X} is

$$D(X_i, X_j) := \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} |x_i - x_j| f_{X_i}(x_i) f_{X_j}(x_j) dx_i dx_j.$$

Then for $X_i \sim \mathcal{N}(\mu_i, \sigma^2)$, $X_j \sim \mathcal{N}(\mu_j, \sigma^2)$, distance between these Normally distributed variables is

$$\begin{aligned} & D_{NN}(X_i, X_j) \\ := & \frac{1}{2\pi\sigma^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} |x_i - x_j| \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma^2}\right) \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma^2}\right) dx_i dx_j \\ = & \frac{2\sigma}{\sqrt{\pi}} \exp\left(-\frac{(\mu_i - \mu_j)^2}{4\sigma^2}\right) + |\mu_i - \mu_j| \operatorname{erf}\left(\frac{|\mu_i - \mu_j|}{2\sigma}\right), \end{aligned} \quad (2.6)$$

where the subscript “NN” in the LHS qualifies the distance as between 2 Normally distributed r.v.s. \square

Proposition 2.1. *SRGG $\mathcal{G}_{m, \mathbf{R}}(\mathbf{V}, \tau)$ drawn with affinity cut-off τ , in the probabilistic metric space \mathcal{X} , with affinity measure $m(G_{ij} | \rho_{ij})$ defined as in Equation 2.3, results from randomly placing Normally distributed random variables with respective means, and a global variance σ^2 , in \mathcal{X} , as long as we set:*

$$\sigma \equiv \frac{1}{\sqrt{2}}; \quad |\mu_i - \mu_j| \equiv |G_{ij} - |\rho_{ij}||.$$

Here matrix $\mathbf{R} = [\rho_{ij}]$.

Proof. We see in Equation 2.6 that setting $2\sigma = \sqrt{2}$, and $|\mu_i - \mu_j| \equiv |G_{ij} - |\rho_{ij}||$, implies distance between Normally distributed r.v.s $X_i \sim \mathcal{N}(\mu_i, \sigma)$, $X_j \sim \mathcal{N}(\mu_j, \sigma)$, with $X_i, X_j \in \mathcal{X}$ is

$$D_{NN}(X_i, X_j) = \frac{\sqrt{2}}{\sqrt{\pi}} \exp\left(-\frac{(G_{ij} - \rho_{ij})^2}{2}\right) + |G_{ij} - \rho_{ij}| \operatorname{erf}\left(\frac{|G_{ij} - |\rho_{ij}||}{\sqrt{2}}\right) = D_{ij},$$

where D_{ij} is defined in Equation 2.5 in terms of the complimentary connection function, (i.e. complimentary affinity measure) of our SRGG, as $D_{ij} := m(G_{ij}\rho_{ij}) + |G_{ij} - |\rho_{ij}||$. \square

Theorem 2.3. *SRGG $\mathcal{G}_{m,\mathbf{R}}(\mathbf{V}, \tau)$, drawn with affinity cut-off τ , in the probabilistic metric space \mathcal{X} , with affinity measure $m(G_{ij}|\rho_{ij})$ defined as in Equation 2.3, is generated by a non-homogeneous Poisson point process.*

Proof. Vertex set of SRGG $\mathcal{G}_{m,\mathbf{R}}(\mathbf{V}, \tau)$ is $\mathbf{V} = \{X_1, \dots, X_p\}$.

For any $i \in \{1, \dots, p\}$, let ball $B(X_i, a)$ be drawn in \mathcal{X} , centred at X_i , with radius a , where $a \in [0, 1]$, given that radius of a ball in \mathcal{X} is a probability.

Let r.v. $N(a) :=$ number of elements of \mathbf{V} inside $B(X_i, a)$.

Then $\forall B(X_i, a) \subset \mathcal{X}$, recalling that $G_{ij} = 0$ if $m(G_{ij}|\rho_{ij}) < \tau$, the expectation of $N(a)$ is:

$$\begin{aligned} \mathbb{E}[N(a)] &= \sum_{j=1}^p f_{X_i}(\mu_i, \sigma) \pi a^2 H(m(G_{ij}|\rho_{ij}) - \tau), & (2.7) \\ &= f_{X_i}(\mu_i, \sigma) \pi a^2 \sum_{j=1}^p H(m(G_{ij}|\rho_{ij}) - \tau), \\ &\equiv f_{X_i}(\mu_i, \sigma) \pi a^2 Q_{\mathbf{R}, \tau} \end{aligned}$$

where the Heaviside function $H(\cdot)$ is defined as

$$\begin{aligned} H(x) &= 1 \quad \text{if } x \geq 0 \\ &= 0 \quad \text{if } x < 0 \end{aligned} \quad (2.8)$$

and we introduce the notation:

$$Q_{\mathbf{R}, \tau} := \sum_{j=1}^p H(m(G_{ij}|\rho_{ij}) - \tau)$$

We further define

$$\lambda_i := f_{X_i}(\mu_i, \sigma) Q_{\mathbf{R}, \tau},$$

where we recall that

$$f_{X_i}(\mu_i, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma^2}\right).$$

Then PP $\{X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_p\}$ approximates a non-homogeneous Poisson Process with intensity λ_i .

□

Remark 2.3. *Thus, our SRGG is underlined by a non-homogeneous Poisson point process that has a location-dependent rate. Hence our SRGG is location-dependent.*

2.5. Dependence of graph on partial correlation matrix \mathbf{R}

SRGG $\mathcal{G}_{m, \mathbf{R}}(\mathbf{V}, \tau)$ is learnt given the partial correlation matrix \mathbf{R} that is itself learnt given the data. In fact, it is the correlation matrix that is learnt given the data, and \mathbf{R} is computed at every update of the correlation matrix. In this section we discuss the learning of the correlation matrix given a multivariate dataset.

Definition 2.4. *Let $\mathbf{X} \in \Xi \subseteq \mathbb{R}^p$ be a p -dimensional observed vector, with $\mathbf{X} = (X_1, \dots, X_p)^T$.*

Let there be n measurements of X_j , $j = 1, \dots, p$, so that the $n \times p$ -dimensional matrix $\mathbf{D} = [x_{ij}]_{i=1; j=1}^{n;p}$ is the data that comprises n measurements of the p -dimensional observable \mathbf{X} .

Let the i -th realisation of \mathbf{X} be \mathbf{x}_i , $i = 1, \dots, n$.

We model \mathbf{X} , so that the set of n realisations of this variable that comprises the data \mathbf{D} , is jointly matrix-Normal, i.e. where this matrix-Normal density is parametrised by

- a mean matrix $\boldsymbol{\mu}^{(n \times p)}$;*
- a covariance matrix $\boldsymbol{\Sigma}_R^{(n \times p)}$, an element of which is the covariance between a pair of rows in \mathbf{D} ;*
- a covariance matrix $\boldsymbol{\Sigma}_C^{(p \times p)}$ that informs on inter-column covariance in data \mathbf{D} .*

We standardise the variable X_j ($j = 1, \dots, p$) by its empirical mean and standard deviation, into Z_j , s.t. the standardised version \mathbf{D}_S of data \mathbf{D} comprises

n measurements of the p -dimensional vector $\mathbf{Z} = (Z_1, \dots, Z_p)^T$. Thus,

$$z_{ij} = \frac{x_{ij} - \bar{x}_j}{\Upsilon_j}, \quad \text{where} \quad \bar{x}_j := \frac{\sum_{i=1}^n x_{ij}}{n}; \Upsilon_j^2 := \frac{\sum_{i=1}^n x_{ij}^2}{n} - \left(\frac{\sum_{i=1}^n x_{ij}}{n} \right)^2.$$

The $n \times p$ -dimensional matrix $\mathbf{D}_S = [z_{ij}]$.

Then we model the joint probability of a set of n measurements of \mathbf{Z} , that comprises the standardised data \mathbf{D}_S , to be matrix-Normal with zero-mean, i.e.

$$\{\mathbf{z}_1, \dots, \mathbf{z}_n\} \sim \mathcal{MN}(\mathbf{0}, \boldsymbol{\Sigma}_R^{(S)}, \boldsymbol{\Sigma}_C^{(S)}),$$

i.e. the likelihood of the covariance matrices $\boldsymbol{\Sigma}_R^{(S)}$ and $\boldsymbol{\Sigma}_C^{(S)}$, given data \mathbf{D}_S , is matrix-Normal:

$$\begin{aligned} \ell(\boldsymbol{\Sigma}_R^{(S)}, \boldsymbol{\Sigma}_C^{(S)} | \mathbf{D}_S) = \\ \frac{1}{(2\pi)^{\frac{np}{2}} |\boldsymbol{\Sigma}_C^{(S)}|^{\frac{p}{2}} |\boldsymbol{\Sigma}_R^{(S)}|^{\frac{n}{2}}} \times \exp \left[-\frac{1}{2} \text{tr} \left\{ (\boldsymbol{\Sigma}_R^{(S)})^{-1} \mathbf{D}_S (\boldsymbol{\Sigma}_C^{(S)})^{-1} (\mathbf{D}_S)^T \right\} \right]. \end{aligned} \quad (2.9)$$

Here $\boldsymbol{\Sigma}_R^{(S)}$ generates the covariance between the standardised variables \mathbf{Z}_i and $\mathbf{Z}_{i'}$, $i, i' = 1, \dots, n$, (while $\boldsymbol{\Sigma}_R$ generates the covariance between \mathbf{X}_i and $\mathbf{X}_{i'}$). Similarly, $\boldsymbol{\Sigma}_C^{(S)}$ generates the correlation between columns of \mathbf{D}_S .

Theorem 2.4. *When the prior on $\boldsymbol{\Sigma}_C^{(S)}$ is uniform, the joint posterior probability density of the correlation matrices $\boldsymbol{\Sigma}_C^{(S)}, \boldsymbol{\Sigma}_R^{(S)}$, given the standardised data \mathbf{D}_S can be marginalised over the $n \times n$ -dimensional inter-row correlation $\boldsymbol{\Sigma}_R^{(S)}$, to yield*

$$\pi(\boldsymbol{\Sigma}_C^{(S)} | \mathbf{D}_S) \propto \frac{1}{c(\boldsymbol{\Sigma}_C^{(S)}) \left| \boldsymbol{\Sigma}_C^{(S)} \right|^{p/2} \left| \mathbf{D}_S (\boldsymbol{\Sigma}_C^{(S)})^{-1} (\mathbf{D}_S)^T \right|^{\frac{n+1}{2}}},$$

where prior on $\boldsymbol{\Sigma}_R^{(S)}$ is the non-informative $\pi_0(\boldsymbol{\Sigma}_R^{(S)}) = \left| \boldsymbol{\Sigma}_R^{(S)} \right|^\alpha$; $\alpha = -\frac{n}{2} - 1$; and $\boldsymbol{\Sigma}_C^{(S)}$ is assumed invertible. Here, $c(\boldsymbol{\Sigma}_C^{(S)})$ is a function of $\boldsymbol{\Sigma}_C^{(S)}$ that normalises the likelihood.

As posterior $\pi(\boldsymbol{\Sigma}_C^{(S)} | \mathbf{D}_S)$ above is obtained for Uniform prior on $\boldsymbol{\Sigma}_C^{(S)}$, the likelihood of $\boldsymbol{\Sigma}_C^{(S)}$ given data \mathbf{D}_S , i.e. pdf of \mathbf{D}_S given $\boldsymbol{\Sigma}_C^{(S)}$ is:

$$\ell(\boldsymbol{\Sigma}_C^{(S)} | \mathbf{D}_S) \equiv f(\mathbf{D}_S | \boldsymbol{\Sigma}_C^{(S)}) \propto \pi(\boldsymbol{\Sigma}_C^{(S)} | \mathbf{D}_S).$$

The proof of this theorem is provided in Section 2 of the attached Supplementary Materials.

Proposition 2.2. *An estimator of the normalisation $c\left(\boldsymbol{\Sigma}_C^{(S)}\right)$ of the posterior $\left[\boldsymbol{\Sigma}_C^{(S)}|\mathbf{D}_S\right]$, given by Theorem 2.4*

$$\hat{c}\left(\boldsymbol{\Sigma}_C^{(S)}\right) = \mathbb{E}_{Z'/_{n/p}} \left[\dots \left[\mathbb{E}_{Z'_{11}} \left[\frac{1}{\left| \left(\mathbf{D}'/\left(\boldsymbol{\Sigma}_C^{(S)}\right)^{-1}\left(\mathbf{D}'\right)^T \right) \right|^{\frac{n'+1}{2}}} \right] \right] \dots \right].$$

Proof. We substitute the sequential computing of expectations with respect to (w.r.t.) distribution of each element of dataset \mathbf{D}' – as suggested in the statement of this Proposition – with computation of the expectation w.r.t. the block \mathbf{D}' of these elements, where \mathbf{D}' abides by the inter-column correlation of $\boldsymbol{\Sigma}_C^{(S)}$. Thus, we approximate the normalisation $c'\left(\boldsymbol{\Sigma}_C^{(S)}\right)$ as:

$$\hat{c}'\left(\boldsymbol{\Sigma}_C^{(S)}\right) = \mathbb{E}_{\mathbf{D}'_S} \left[\frac{1}{\left| \left(\mathbf{D}'/\left(\boldsymbol{\Sigma}_C^{(S)}\right)^{-1}\left(\mathbf{D}'\right)^T \right) \right|^{\frac{n'+1}{2}}} \right].$$

We consider the sample of k number of $n' \times p$ -dimensional data sets $\{\mathbf{D}'_1, \dots, \mathbf{D}'_k\}$, where \mathbf{D}'_q abides by inter-column correlation $\boldsymbol{\Sigma}_C^{(S)} \forall q = 1, \dots, k$, s.t. $\mathbf{D}'_q\left(\boldsymbol{\Sigma}_C^{(S)}\right)^{-1}\left(\mathbf{D}'_q\right)^T$ is positive definite $\forall q = 1, \dots, k$, at each t .

Then an estimator of $\hat{c}'\left(\boldsymbol{\Sigma}_C^{(S)}\right)$ is

$$\hat{c}_C^{(S)} := \frac{1}{K} \sum_{k=1}^K \frac{1}{\left| \left(\mathbf{D}'_k/\left(\boldsymbol{\Sigma}_C^{(S)}\right)^{-1}\left(\mathbf{D}'_k\right)^T \right) \right|^{\frac{n'+1}{2}}}. \quad (2.10)$$

□

Generation of a randomly sampled $n' \times p$ -sized data set \mathbf{D}'_k , with column correlation $\boldsymbol{\Sigma}_C^{(S)}$, is undertaken.

3. Bayesian Inference on SRGG and Correlation Matrix using MCMC, and the Compounding of Processes Underlying the SRGG

We perform Bayesian inference on the matrix $\mathbf{R} = [\rho_{ij}]$ of partial correlations between each distinct pair of the observed variables in data \mathbf{D}_S , and simultaneously, on the \mathbf{R} -dependent SRGG $\mathcal{G}_{m,\mathbf{R}}(\mathbf{V}, \tau)$ drawn in probabilistic metric

space \mathcal{X} with affinity measure $m(G_{ij}|\rho_{ij})$, on vertex set \mathbf{V} , with cutoff probability τ . The inference that we undertake is Markov Chain Monte Carlo-based, i.e. MCMC-based. In particular it is an implementation of the Metropolis-with-2-block-update.

Remark 3.1. *In the implementation of the Metropolis-with-2-block-update, the inter-column correlation matrix $\Sigma_C^{(C)}$ of the data is first updated – at which the updated partial correlation matrix \mathbf{R} is computed – for the SRGG to be then updated, at this updated \mathbf{R} .*

Remark 3.2. *Posterior probability $\pi(\Sigma_C^{(S)}|\mathbf{D}_S)$ of inter-column correlation matrix $\Sigma_C^{(S)}$, given data \mathbf{D}_S , as given in Theorem 2.4, implies that the sequence of realisations of $\Sigma_C^{(S)}$, at successive iterations of the MCMC chain:*

$$\{\Sigma_C^{(S)0}, \Sigma_C^{(S)1}, \Sigma_C^{(S)2}, \dots\},$$

is a continuous-valued, discrete time stochastic process, with underlying probability density $\pi(\Sigma_C^{(S)}|\mathbf{D}_S)$.

Definition 3.1. *Given a learnt value of the inter-column correlation matrix $\Sigma_C^{(S)}$, to compute the value ρ_{ij} of the partial correlation R_{ij} between r.v.s X_i and X_j , we first invert the $p \times p$ -dimensional matrix $\Sigma_C^{(S)}$ to yield the precision matrix:*

$$\Psi := \left(\Sigma_C^{(S)}\right)^{-1}; \Psi = [\psi_{ij}],$$

s.t. the partial correlation matrix $\mathbf{R} = [\rho_{ij}]$, where ρ_{ij} is

$$\rho_{ij} = -\frac{\psi_{ij}}{\sqrt{\psi_{ii}\psi_{jj}}}, \quad i \neq j, \quad (3.1)$$

and $\rho_{ii} = 1$ for $i = j$.

Proposition 3.1. *Following Proposition 2.1, the non-homogeneous Poisson point process with \mathbf{R} -dependent intensities, that generates the SRGG $\mathcal{G}_{m,\mathbf{R}}(\mathbf{V}, \tau)$, is compounded with the continuous-valued stochastic process $\{\Sigma_C^{(S)t}\}_{t \in \{0,1,2,\dots\}}$ (discussed in remark 3.2), with underlying density $\pi(\Sigma_C^{(S)}|\mathbf{D}_S)$ that is defined in Theorem 2.4.*

3.1. MCMC-driven definition of the 95% HPD credible regions on the learnt SRGG

To acknowledge uncertainties in the Bayesian learning of the sought SRGG, where the uncertainties can be identified with a Bayesian 95% HPD credible region, we can suggest that the learnt SRGG only contains those edges, the affinity measures of which exceed a global probability of 0.05. Then we are basically defining $\tau = 0.05$ to define the edge set of the sought SRGG (see Definition 2.3).

Within our MCMC-based inference, the random graph $\mathcal{G}_{m, \mathbf{R}^{(t)}}(\mathbf{V}, \tau)$, is sampled in the t -th iteration of the inference scheme, given the inter-column correlation matrix $\Sigma_C^{(S)t}$ updated in this iteration, using which partial correlation matrix $\mathbf{R}^{(t)} = [\rho_{ij}^{(t)}]$ is updated in the t -th iteration; $t = 0, 1, \dots, N_{iter}$.

In the t -th iteration, let affinity between r.v.s X_i and X_j ($X_i \neq X_j; X_i, X_j \in \mathbf{V}$) be given by the marginal $m(G_{ij} = g_{ij}^{(t)} | \rho_{ij}^{(t)})$.

Then the graphical model of data \mathbf{D}_S is defined as the graph on vertex set \mathbf{V} , that includes edge between X_i and X_j , ($X_i \neq X_j; X_i, X_j \in \mathbf{V}$), if sample estimate of $\mathbb{E}[m(G_{ij} | \rho_{ij}^{(t)})]$ exceeds $\tau = 0.05$. We delineate this definition formally in Definition 3.2.

A sample estimate of $\mathbb{E}[m(G_{ij} | \rho_{ij}^{(t)})]$ is

$$\hat{m}(G_{ij} | \rho_{ij}^{(t)}) := N_{ij},$$

where N_{ij} is the the fraction of the post-burnin (N_{post}) number of iterations (where $N_{post} < N_{iter} + 1$), in which the ij -th edge exists, i.e. in which G_{ij} takes the value 1, $\forall i, j = 1, 2, \dots, p$, $i \neq j$. Thus, variable N_{ij} takes the value

$$n_{ij} := \frac{\sum_{t=N-N_{post}+1}^N g_{ij}^{(t)}}{N_{post}}, \quad i < j; i, j = 1, \dots, p. \quad (3.2)$$

Remark 3.3. N_{ij} is the (post-burnin) sample mean of the affinity measure between X_i and X_j ($X_i, X_j \in \mathbf{V}$), i.e. of the marginal of edge G_{ij} , conditional on the learnt inter-column (partial) correlation matrix of the given dataset.

We define the $p \times p$ “edge-parameter matrix” as

$$\mathcal{N} := [n_{ij}].$$

Indeed, the N_{ij} parameter is a function of the partial correlation ρ_{ij} , that is itself learnt given this data, but for the sake of notational brevity, we do not include this explicit \mathbf{R} dependence in our notation.

Definition 3.2. The set of SRGGs:

$$\{\mathcal{G}_{m, \mathbf{R}^{(t)}}(\mathbf{V}, \tau)\}_{t=N-N_{post}+1}^N$$

sampled during the post-burnin part of any MCMC chain, on the vertex set $\mathbf{V} = \{X_1, \dots, X_p\}$, for $\tau = 0.05$, defines the graphical model $\hat{G}_{\mathcal{N}, \mathbf{D}_S}(\mathbf{V}, 0.05)$ of the data \mathbf{D}_S , learnt within 95% HPD credible region, as the graph on vertex set \mathbf{V} , with the edge set:

$$\hat{E}_p = \{G_{ij} : n_{ij} \geq 0.05; X_i \neq X_j; X, X_j \in \mathbf{V}\}.$$

Here, the (post-burnin) sample mean of the affinity measure between X_i and X_j is n_{ij} s.t. edge-parameter matrix is $\mathcal{N} = [n_{ij}]$.

3.2. Computational Details of Metropolis-with-2-block-update

Theorem 2.4 gives the posterior probability density of correlation matrix $\Sigma_C^{(S)}$, given data \mathbf{D}_S . In our Metropolis-with-2-block-update based inference, we update $\Sigma_C^{(S)}$ —at which the partial correlation matrix \mathbf{R} is computed. Given this updated \mathbf{R} , we then update the graph.

In our learning of the $p \times p$ -dimensional inter-column correlation matrix $\Sigma_C^{(S)} = [s_{ij}]$ —where s_{ij} is the value of r.v. S_{ij} —the $\frac{p^2 - p}{2}$ non-diagonal elements of the upper (or lower) triangle are learnt, i.e. the parameters $S_{12}, S_{13}, \dots, S_{1p}, S_{23}, \dots, S_{p-1,p}$ are learnt.

In the t -th iteration, S_{ij} is proposed from a Truncated Normal density that is left truncated at -1 and right truncated at 1, as

$$S_{ij} = s_{ij}^{(t^*)} \sim \mathcal{TN}(s_{ij}^{(t^*)}; s_{ij}^{(t-1)}, \sigma_{ij}^2, -1, 1), \quad \forall i, j = 1, \dots, p; i \neq j,$$

where $\sigma_{ij} = \sigma_0^2 \forall i, j$, is the experimentally chosen variance, and the proposal mean is the current value $s_{ij}^{(t-1)}$ of S_{ij} at the end of the $t - 1$ -th iteration.

At the 2nd block of the t -th iteration, the SRGG is updated, given the recently updated partial correlation matrix $\mathbf{R}^{(t)} = [\rho_{ij}^{(t)}]$, s.t. the proposed edge variable connecting the i -th to the j -th vertex is

$$G_{ij} = g_{ij}^{(t^*)} \sim \text{Bernoulli}(g_{ij}^{(t^*)}; \rho_{ij}^{(t)}).$$

This update also involves the likelihood of the edge parameters G_{ij} and variance parameters v_{ij} , $\forall i \neq j; i, j \in \{1, \dots, p\}$, given the updated partial correlation matrix. We recall from Equation 2.1 that this likelihood is

$$\begin{aligned} \ell(G_{12}, \dots, G_{1p}, G_{23}, \dots, G_{2p}, \dots, G_{p,p-1}, v_{12}, \dots, v_{1p}, v_{23}, \dots, v_{2p}, \dots, v_{p,p-1} | \mathbf{R}) \\ = \prod_{i \neq j; i, j=1}^p \frac{1}{\sqrt{2\pi v_{ij}}} \exp \left[-\frac{(G_{ij} - |\rho_{ij}^{(t)}|)^2}{2v_{ij}} \right]. \end{aligned}$$

The ij -th variance parameter is assigned a proposed value of

$$v_{ij}^{(t\star)} \sim \mathcal{N}(v_{ij}^{(t\star)}; v_{ij}^{(t-1)}, w_{ij}^2),$$

where w_{ij}^2 is the experimentally chosen variance, and the mean is the current value of the v_{ij} variable.

As suggested in Theorem 2.4, the correlation learning involves computing $(\Sigma_C^{(S)})^{-1}$, $|\Sigma_C^{(S)}|$ and $|\mathbf{D}_S (\Sigma_C^{(S)})^{-1} (\mathbf{D}_S)^T|$, in every iteration. This calls for Cholesky decomposition of $\Sigma_C^{(S)}$ as $\mathbf{L}_C^{(S)} (\mathbf{L}_C^{(S)})^T$, and of $\mathbf{D}_S (\Sigma_C^{(S)})^{-1} (\mathbf{D}_S)^T$, into the (lower) triangular matrix \mathbf{L} and \mathbf{L}^T , while implementing ridge adjustment (Wothke, 1993). The latter computation follows the inversion of $\Sigma_C^{(S)}$ into $(\Sigma_C^{(S)})^{-1}$, which is undertaken using a forward substitution algorithm. (see Section 10 of the Supplementary Materials).

Once the set of SRGGs sampled during the post-burnin part of the MCMC chain are identified ($\{\hat{\mathcal{G}}_{m, \mathbf{R}^{(t)}}(\mathbf{V}, \tau)\}_{t=N-N_{post}+1}^N$), graphical model $\hat{\mathcal{G}}_{\mathcal{N}, \mathbf{D}_S}(\mathbf{V}, 0.05)$ of the data \mathbf{D}_S is then constructed, using this set of SRGGs.

4. Inter-graph distance metric

We compute the Hellinger distance between the posterior probability of the graphical model $\hat{\mathcal{G}}_{\mathcal{N}_1, \mathbf{D}_1}(\mathbf{V}, 0.05)$ of the data \mathbf{D}_1 , learnt within 95% HPD credible region, and the posterior probability of the similarly learnt graphical model $\hat{\mathcal{G}}_{\mathcal{N}_2, \mathbf{D}_2}(\mathbf{V}, 0.05)$ of the data \mathbf{D}_2 . Distance between the pair of uncertainty-included graphical models is computed using the Hellinger metric, normalised by the uncertainty in the learning of each graphical model, where such uncertainty is defined below (Definition 4.3). Here data \mathbf{D}_1 comprises n_1 rows and p columns, while data \mathbf{D}_2 comprises n_2 rows and p columns. (As we soon explain, we compute the Hellinger distance between the marginal posterior probability of the edges in each of the considered pair of graphical models).

The inter-graph distance is computed, to inform on the absolute of the correlation between the multivariate, disparately-sized datasets \mathbf{D}_1 and \mathbf{D}_2 ; in effect, the exercise can address the possible independence of the *pdfs* that the two datasets are sampled from. This is of course a hard question to address when the data comprise measurements of a high-dimensional vector-valued observable.

Definition 4.1. *Square of Hellinger distance between two probability density functions $g(\cdot)$ and $h(\cdot)$ over a common domain $\mathcal{X} \in \mathbb{R}^m$, with respect to a*

chosen measure, is

$$\begin{aligned}
D_H^2(g, f) &= \int \left(\sqrt{g(\mathbf{x})} - \sqrt{h(\mathbf{x})} \right)^2 d\mathbf{x} \\
&= \int g(\mathbf{x}) d\mathbf{x} + \int h(\mathbf{x}) d\mathbf{x} - 2 \int \sqrt{g(\mathbf{x})} \sqrt{h(\mathbf{x})} d\mathbf{x} \\
&= 2 \left(1 - \int \sqrt{g(\mathbf{x})} \sqrt{h(\mathbf{x})} d\mathbf{x} \right). \tag{4.1}
\end{aligned}$$

The Hellinger distance is closely related to the Bhattacharyya distance (Bhattacharyya, 1943) between two densities: $D_B(g, f) = -\log \left[\int \left(\sqrt{g(\mathbf{x})} \sqrt{h(\mathbf{x})} \right)^2 d\mathbf{x} \right]$.

Definition 4.2. Consider the marginal posterior probability density of all the graph edge parameters $\{G_{ij}\}_{i \neq j; i, j=1}^p$ given the partial correlation matrix \mathbf{R}_q (that is itself updated given the data $\mathbf{D}_S^{(q)}$); $q = 1, 2$.

In the t -th iteration, value of the marginal posterior of all the edges $\{G_{ij}^{(qt)}\}_{i \neq j; i, j=1}^p$, in the q -th SRGG, given \mathbf{R}_{qt} , is:

$$m(G_{11}^{(qt)}, G_{12}^{(qt)}, \dots, G_{p-p-1}^{(qt)} | \mathbf{R}_{qt}), \quad t = 0, \dots, N_{iter}.$$

Given the availability of the value of this marginal posterior density, only at discretely sampled points in its support, (sampled at discrete times t), the integral in the definition of the Hellinger distance is replaced by a sum in our computation of the distance.

Then for the q -th dataset, the marginal posterior of all graph edge parameters in the t -th iteration is:

$$u_q^{(t)} := m(G_{11}^{(qt)}, G_{12}^{(qt)}, \dots, G_{p-p-1}^{(qt)} | \mathbf{R}_{qt}),$$

which is employed to compute square of the (discretised version of the) Hellinger distance between the two datasets as

$$D_H^2(u_1, u_2) = \frac{\sum_{t=N_{burnin}+1}^{N_{iter}} \left(\sqrt{u_1^{(t)}} - \sqrt{u_2^{(t)}} \right)^2}{N_{iter} - N_{burnin}}, \tag{4.2}$$

The Bhattacharyya distance can be similarly discretised.

However, MCMC does not provide normalised posterior probability densities – we may employ Uniform (over identified finite intervals) priors on the variance parameters, the marginalised posterior probability of the edge parameters is known only up to an unknown scale.

Remark 4.1. In the t -th iteration, MCMC provides value of logarithm $\ln(u_q^{(t)})$ of the un-normalised posterior of the edges of the graph given the q -th data ($q = 1, 2$). Hence the Hellinger distance between the 2 datasets that we compute is only known upto a constant normalisation S that we use to scale both $u_1^{(t)}$ and $u_2^{(t)}$, $\forall t = 0, \dots, N_{iter}$.

Proposition 4.1. Unknown normalisation S that normalises $u_1^{(t)}$ and $u_2^{(t)}$, is chosen to ensure that the scaled, log marginal of all graph edges in the t -th iteration, is ≤ 0 , s.t. $\exp\left(\frac{\ln(u_m^{(t)})}{s}\right) \in (0, 1]$. Therefore we choose the global scale S as:

$$s := \max\{\ln(u_1^{(0)}), \ln(u_1^{(1)}), \dots, \ln(u_1^{(N_{iter})}), \ln(u_2^{(0)}), \dots, \ln(u_2^{(N_{iter})})\}. \quad (4.3)$$

Remark 4.2. Squared Hellinger distance $D_H^2(u_1, u_2)$ between discretised posterior probability densities of 2 graphical models, computed using $\exp(\ln(u_q^{(t)})/s)$ in Equation 4.2, is affected by scaling parameter S . This scale dependence is mitigated in our definition of the distance between 2 graphical models as the difference between the ratio of this computed $D_H(u_1, u_2)$, to the scaled uncertainty inherent in one graphical model, and the ratio of $D_H(u_1, u_2)$, to the scaled uncertainty in the other learnt graphical model. Such scaled uncertainty in a learnt graphical model is defined in Definition 4.3.

Definition 4.3. The scaled (by a scale parameter $S = s$) uncertainty in learnt graphical model $\hat{G}_{\mathcal{N}_q, \mathbf{D}_q}(\mathbf{V}, 0.05)$ of data set \mathbf{D}_q , with edge-parameter matrix \mathcal{N}_q , is defined as

$$D_{max,s}(q) := \max\{\exp(\ln(u_q^{(0)})/s), \exp(\ln(u_q^{(1)})/s), \dots, \exp(\ln(u_q^{(N_{iter})})/s)\} - \min\{\exp(\ln(u_q^{(0)})/s), \exp(\ln(u_q^{(1)})/s), \dots, \exp(\ln(u_q^{(N_{iter})})/s)\}, \quad (4.4)$$

Thus, $D_{max,s}(q)$ provides separation between the maximal and minimal (scaled values of) posteriors of graphs, generated in the MCMC chain run using the q -th dataset. Therefore $D_{max,s}(q)$ defines uncertainty of the graphical model learnt for this dataset.

Definition 4.4. For edge-parameter matrix \mathcal{N}_q , for dataset \mathbf{D}_q , $q = 1, 2$, separation between the two corresponding graphical models on vertex set \mathbf{V} , learnt with uncertainty defined as in Definition 4.3 is

$$\begin{aligned} & \delta(\hat{\mathcal{G}}_{\mathcal{N}_1, \mathbf{D}_1}(\mathbf{V}, 0.05), \hat{\mathcal{G}}_{\mathcal{N}_2, \mathbf{D}_2}(\mathbf{V}, 0.05)) \\ & := \left| \sqrt{D_H^2(u_1, u_2)/D_{max,s}(1)} - \sqrt{D_H^2(u_1, u_2)/D_{max,s}(2)} \right| \\ & = D_H(u_1, u_2) \left| \frac{1}{D_{max,s}(1)} - \frac{1}{D_{max,s}(2)} \right|, \end{aligned} \quad (4.5)$$

where the Hellinger distance $D_H(u_1, u_2)$, between the 2 graphical models, is defined in Equation 4.2 and $D_{max,s}(q)$ is the uncertainty in the graphical model for data \mathbf{D}_q , as defined in Equation 4.4, computed at the chosen value s of scale S (defined in Equation 4.3).

Alternatively, we could define a (discretised version of the) odds ratio of unscaled logarithm of the unnormalised posterior densities of the graphical models learnt using MCMC, given the two datasets, as $\int (\log(g(\mathbf{x})) - \log(h(\mathbf{x}))) d\mathbf{x}$; such is then a divergence measure that we define as

$$O_\pi(u_1, u_2) := \sum_{t=N_{burnin}+1}^{N_{iter}} [\log(u_1^{(t)}) - \log(u_2^{(t)})]. \quad (4.6)$$

4.1. Suggested inter-graph separation $\delta(\cdot, \cdot)$, is an inter-graph distance

Theorem 4.1. Let $\delta(\hat{\mathcal{G}}_{\mathcal{N}_1, \mathbf{D}_1}(\mathbf{V}, 0.05), \hat{\mathcal{G}}_{\mathcal{N}_2, \mathbf{D}_2}(\mathbf{V}, 0.05))$ be the separation defined as in Equation 4.5, between 2 uncertainty-included graphical models, defined over vertex set \mathbf{V} , learnt for datasets \mathbf{D}_1 and \mathbf{D}_2 . Here the graphical model $\hat{\mathcal{G}}_{\mathcal{N}_q, \mathbf{D}_q}(\mathbf{V}, 0.05)$ is an element of space Ω , $q = 1, 2$.

Then our definition of this inter-graph separation $\delta : \Omega \times \Omega \rightarrow \mathbb{R}_{\geq 0}$, is a distance function, or a metric.

The proof of this theorem is provided in Section 3 of the attached Supplementary Materials.

4.2. Absolute correlation between 2 multivariate datasets, from distance between their graphical models

In this section, we introduce a model for the absolute correlation between 2 multivariate datasets, for which the uncertainty-included graphical models are learnt, allowing for the inter-graph distance $\delta(\cdot, \cdot)$ to be computed.

Proposition 4.2. *For a given value of the inter-graph distance $\delta(\mathcal{G}_{u,1}, \mathcal{G}_{u,2}) \in [0, \infty)$, between 2 learnt graphical models $\mathcal{G}_{u,2}, \mathcal{G}_{u,1} \in \Omega$, defined over vertex set $\{1, \dots, p\}$, where the graphical model $\mathcal{G}_{u,\cdot}$ is learnt given data \mathbf{D} , a model for the absolute value of the correlation $|\text{corr}(\mathbf{Z}_1, \mathbf{Z}_2)|$ between the p -dimensional vector-valued observable \mathbf{Z}_1 , (n_1 measurements of which comprise dataset indexed by 1), and the p -dimensional observable \mathbf{Z}_2 , (n_2 measurements of which comprise dataset indexed by 2), is*

$$\delta(\mathcal{G}_{u,1}, \mathcal{G}_{u,2}) = -\log(|\text{corr}(\mathbf{Z}_1, \mathbf{Z}_2)|),$$

$$\text{s.t. } |\text{corr}(\mathbf{Z}_1, \mathbf{Z}_2)| = \exp[-\delta(\mathcal{G}_{u,1}, \mathcal{G}_{u,2})] \in (0, 1].$$

5. Changes undertaken to facilitate the learning of large networks

When our interest is in learning a graphical model on a vertex set of cardinality $p \gtrsim 20$, it implies that such learning, if it is to be undertaken according to the methodology described in the previous section, will demand MCMC-based inference on the $\gtrsim 200$ distinct off-diagonal elements of the correlation matrix $\Sigma_C^{(S)} = [s_{ij}]$ (where S_{ij} represents correlation between the X_i and X_j variables in the dataset); MCMC-based learning of more than about 200 parameters is difficult. Again, for $p \gtrsim 500$, Cholesky decomposition of the $p \times p$ -dimensional inter-column correlation matrix $\Sigma_C^{(S)}$, (leading to its inversion for example) is not easy (i.e. it is a challenge to achieve numerical robustness as the matrix dimensionality exceeds about 500×500). This renders computation of the likelihood in Theorem 2.4 difficult, and the numerical computation of the precision matrix $(\Sigma_C^{(S)})^{-1} = \Psi = [\psi_{ij}]$ is also difficult for $p \gtrsim 500$, where ψ_{ij} is employed to compute the partial correlation matrix \mathbf{R} according to Equation 3.1.

Remark 5.1. *When learning a network with $\gtrsim 500$ vertices as an SRGG drawn in a probabilistic metric space, on vertex set \mathbf{V} with cardinality p , with a cut-off on the affinity measure (\equiv edge marginals) of τ , we learn the SRGG given the correlation matrix $\Sigma_C^{(S)}$ than the partial correlation matrix \mathbf{R} (of the given*

data \mathbf{D}_S that hosts n standardised measurements of each of the $p \gtrsim 500$ r.v.s), since it is hard to compute inverse $(\Sigma_C^{(S)})^{-1}$ of the large $(\Sigma_C^{(S)})^{(p \times p)}$, to thereby compute \mathbf{R} .

Thus, the network is learnt as the SRGG $\mathcal{G}_{m, \Sigma_C^{(S)}}(\mathbf{V}, \tau)$.

Remark 5.2. When learning a network with $\gtrsim 500$ vertices given data \mathbf{D}_S that hosts n standardised measurements of each of $p \gtrsim 500$ r.v.s X_1, \dots, X_p , we — eschew MCMC-based inference on the large $(\Sigma_C^{(S)})^{(p \times p)}$ inter-column correlation matrix, and

— employ empirical estimate of s_{ij} instead, where $\Sigma_C^{(S)} = [s_{ij}]$ with $s_{ij} := \frac{\sum_{k=1}^n x_{ik}x_{jk}}{n} - \frac{\sum_{k=1}^n x_{ik}}{n} \frac{\sum_{k=1}^n x_{jk}}{n}$. Here, k -th measured value of X_i is x_{ik} , $i = 1, \dots, p$; $k = 1, \dots, n$.

Hence in the notation of the network learnt as SRGG $\mathcal{G}_{m, \Sigma_C^{(S)}}(\mathbf{V}, \tau)$, correlation matrix has no dependence on any iteration index.

Remark 5.3. When learning the graphical model of a given dataset \mathbf{D}_S that hosts n standardised measurements of each of $p \gtrsim 500$ r.v.s X_1, \dots, X_p , we eschew MCMC-based inference on an SRGG in every iteration. The sought graphical model is learnt as the network $\mathcal{G}_{m, \Sigma_C^{(S)}}(\mathbf{V}, \tau)$ which is itself an SRGG with connection function, or the affinity measure between the i -th and j -th nodes, given by the marginal posterior of G_{ij} , given correlation $S_{ij} = s_{ij}$ between X_i and X_j , i.e. by:

$$m(G_{ij}|S_{ij}).$$

Indeed, as the MCMC-based inference is not relevant any more, there is only a single value of the marginal posterior $m(G_{ij}|S_{ij})$ of the edge parameter G_{ij} , between the i -th and j -th nodes, (given the correlation S_{ij}). So we do not require to define the connection function in terms of (a sample estimate of) the expected value of the marginal.

Remark 5.4. Graphical model of dataset with inter-column correlation matrix Σ_C , on vertex set \mathbf{V} , with cutoff probability τ , is learnt as the network $\mathcal{G}_{m, \Sigma_C}(\mathbf{V}, \tau)$ with one single identified connection function or affinity function $m(\cdot)$.

Thus, we learn a network as an SRGG without uncertainties.

5.1. Inter-network distance

However, in Definition 4.4, distance between graphical models learnt of a pair of datasets, is defined as the Hellinger distance normalised by the uncertainty in the learning of each graphical model. So in the absence of uncertainty in learning the network as an SRGG, how can we define an inter-network distance? In fact, the very discretised representation of the Hellinger distance between the marginal posteriors of the two graphs, over the MCMC iterations, (see Equation 4.2), stands challenged, when only one marginal value of the SRGG is computed for each given dataset.

Proposition 5.1. *For vertex set $\mathbf{V} = \{X_1, \dots, X_p\}$, distance $\Delta(\cdot, \cdot)$ between network $\mathcal{G}_{m, \Sigma_C^{(1)}}(\mathbf{V}, \tau)$. given a dataset with inter-column correlation matrix $\Sigma_C^{(1)}$, and the network $\mathcal{G}_{m, \Sigma_C^{(2)}}(\mathbf{V}, \tau)$ learnt given dataset with inter-column correlation matrix $\Sigma_C^{(2)}$, is defined as the (discretised) Hellinger distance between the edge marginals of each network, given the respective inter-node correlation structure, i.e. as*

$$\Delta(\mathcal{G}_{m, \Sigma_C^{(1)}}(\mathbf{V}, \tau), \mathcal{G}_{m, \Sigma_C^{(2)}}(\mathbf{V}, \tau)) := D_H(u_1, u_2),$$

where for the q -th dataset, ($q = 1, 2$), the marginal posterior of the ij -th edge parameter $G_{ij}^{(q)}$ given the ij -th correlation parameter $S_{ij}^{(q)}$ is $m(G_{ij}^{(q)} | S_{ij}^{(q)})$, for $i > j; i, j \in \{1, 2, \dots, p\}$, s.t.

$$u_q^{(ij)} := m(G_{ij}^{(q)} | S_{ij}^{(q)}),$$

which is employed to compute square of the (discretised version of the) Hellinger distance between the two datasets as

$$D_H^2(u_1, u_2) = \frac{\sum_{i=1}^{p-1} \sum_{j=i+1}^p \left(\sqrt{u_1^{(ij)}} - \sqrt{u_2^{(ij)}} \right)^2}{p(p-1)/2}, \quad (5.1)$$

The cut-off probability on this marginal posterior is τ in the network learnt as SRGG $\mathcal{G}_{m, \Sigma_C^{(s)}}(\mathbf{V}, \tau)$. Depending on the network at hand, we may decide on the value of τ ; for example, in the human disease-disease network that we learn in Section 8, we produce the network using $\tau = 0.1$.

6. Implementation on real data

In this section we discuss applications of our method to datasets on 12 vino-chemical attributes of two samples of 1599 red and 4898 white wines, grown in the Minho region of Portugal (referred to a “vinho verde”); these data have been considered by Cortez et al. (1998) and discussed in <https://onlinecourses.science.psu.edu/stat857/node/223> (hereon PSU). Each dataset consists of 12 columns that bear information on attributes that are assigned the following names: “fixed acidity” (X_1), “volatile acidity” (X_2), “citric acid” (X_3), “residual sugar” (X_4), “chlorides” (X_5), “free sulphur dioxide” (X_6), “total sulphur dioxide” (X_7), “density” (X_8), “pH” (X_9), “sulphates” (X_{10}), “alcohol” (X_{11}) and “quality” (X_{12}). Then the n -th row and i -th column of the data matrix carries measured/assigned value of the i -th property of the n -th wine in the sample, where $i = 1, \dots, 12$ and $n = 1, \dots, n_{orig} = 1599$ for the red wine data $\mathbf{D}_{orig}^{(red)}$, while $n = 1, \dots, n_{orig} = 4898$ for the white wine data $\mathbf{D}_{orig}^{(white)}$. We refer to the i -th vinous property to be X_i . Then $X_i \in \mathbb{R}_{\geq 0} \forall i = 1, \dots, 11$, while X_{12} that denotes the perceived “quality” of the wine is a categorical variable. Each wine in these samples was assessed by at least three experts who graded the wine on a categorical scale of 0 to 10, in increasing order of excellence. The resulting “sensory score” or value of the “quality” parameter was a median of the expert assessments (Cortez et al., 1998). We seek the graphical model given each of the wine data sets, in which the relationship between any X_i and X_j is embodied, $i \neq j$; $i, j = 1, \dots, 12$. Thus, we seek to find out how the different vino-chemical attributes affect each other, as well as the quality of the wine, in the data at hand. Here, X_1, \dots, X_{11} are real-valued, while X_{12} is a categorical variable, and our methodology allows for the learning of the graphical model of a data set that in its raw state bears measurements of variables of different types. In fact, we standardise our data, s.t. X_i is standardised to Z_i , $i = 1, \dots, p$, $p = 12$. We work with only a subset data set, (comprising only $n < n_{orig}$ rows of the available $\mathbf{D}_{orig}^{(i)}$; $n = 300$ typically). Thus, the data sets with n rows, containing Z_i values, ($i = 1, \dots, p = 12$), are $n \times p$ -dimensional matrices each; we refer to these data sets that we work with, as $\mathbf{D}_S^{(white)}$ and $\mathbf{D}_S^{(red)}$, respectively for the white and red wines. Our aim is to learn the between-column correlation matrix $\Sigma_S^{(m)}$ given data $\mathbf{D}_S^{(m)}$, and simultaneously learn the graphical model of this data using MCMC-based inference within the methodology that we discuss above, to then compute the inter-graph distance, and the inter-data correlation thereafter; $m = white, red$.

The motivation behind choosing these data sets are basically three-fold. Firstly, we sought multivariate, rectangularly-shaped, real-life data, that would admit graphical modelling of the correlations between the different variables in the data. Also, we wanted to work with data, results from – at least a part of – which exists in the literature. Comparison of these published results, with our independent results then illustrates strengths of our method. Thirdly, treating the red and white wine data as data realised at different experimental conditions, we would want to address the question of the distance between these data, and

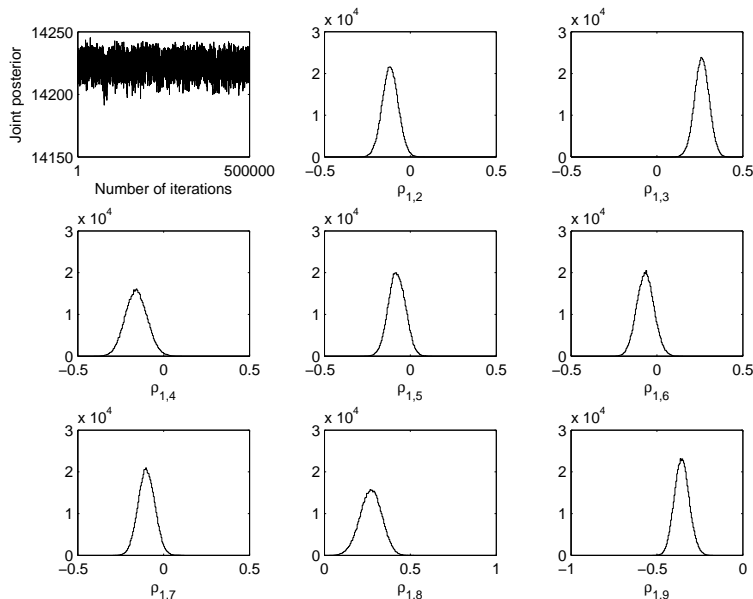


FIG 2. Figure to demonstrate convergence of the MCMC chain that we run with the white wine data. Top left panel: trace of the joint posterior probability density of the elements of the upper triangle of the between-columns correlation matrix of the standardised version of the real data $\mathbf{D}_S^{(white)}$ on Portuguese white wine samples (Cortez et al., 1998); this data has $n = 300$ rows and $p = 12$ columns, and is constructed as a randomly sampled subset of the original data, the sample size of which is 4898. All other panels: histogram representations of marginal posterior probability densities of some of the partial correlation parameters computed using the correlation matrix learnt given data $\mathbf{D}_S^{(white)}$.

we propose to do this by computing the distance between the graphical models of the two data sets. Hence our choice of the popular Portuguese red and white wine data sets, as the data that we implement to illustrate our method on. It is to be noted that a rigorous vinaceous implications of the results, is outside the scope and intent of this paper. However, we will make a comparison of our results with the results of the analysis of white wine data that is reported in PSU, though literature precludes analysis of the red wine data.

6.1. Results given data $\mathbf{D}_S^{(white)}$

Figure 2 presents results that demonstrate convergence of the MCMC chain run to learn the correlation structure and graphical model given the white wine data $\mathbf{D}_S^{(white)}$. Its top left-hand panel displays trace of the joint posterior probability density of all learnt inter-column correlation (S_{ij}) parameters of this data while marginal marginal posterior probabilities of some of the partial correlation

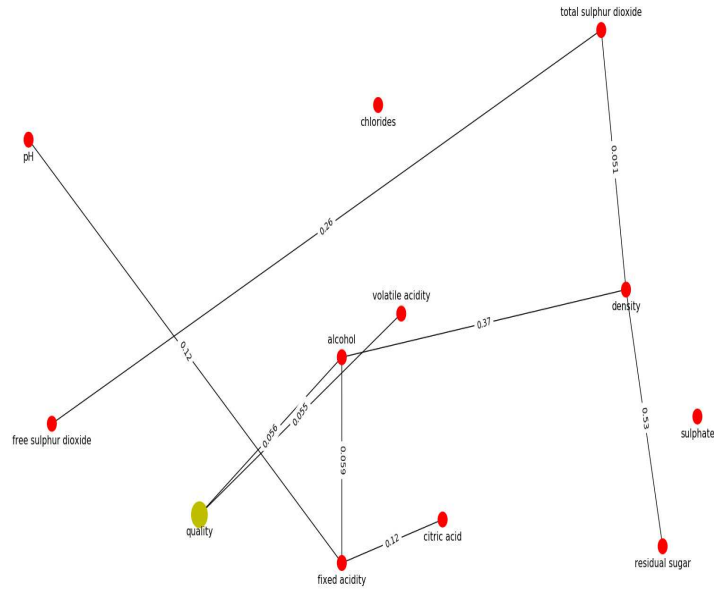


FIG 3. Figure showing graphical model of $\mathbf{D}_S^{(white)}$, of the real data on Portuguese white wine samples (Cortez et al., 1998). The nodes have been placed randomly inside a box. Each of the first 11 columns of this data gives the measured value of each of 11 different vino-chemical properties of the wines in the sample – marked as nodes in the graph above, by filled red (or grey in the printed version) circles, with the name of the property included in the vicinity of the respective node. The 12-th column in the data includes values of the assessed quality of a wine in the sample, (a node that we mark with a green circle in the electronic version; the bigger grey circle in a monochromatic version of the paper). The estimate of the probability for an edge to exist in the post-burnin sample of graphs generated in our MCMC-based inferential scheme, is marked against an existing edge, where edges with such probabilities that are < 0.05 are omitted from this graphical model, (see Section 3.1).

parameters, are presented as histograms in the other panels. Having learnt the correlation structure, the SRGG given the partial correlation matrix updated in an iteration, is then learnt. Figure 9 in Supplementary Materials presents traces of the of some of the edge (G_{ij}) and variance (v_{ij}) parameters of this SRGG learning. Then at the end of the chain, using the learnt SRGGs, graphical model of this data is constructed; this is presented in Figure 3.

6.1.1. Comparing against earlier work done with white wine data

The graphical model of the white wine data presented in Fig 3, is strongly corroborated by the simple empirical correlations between pairs of different vino-chemical properties that is noticed in the “scatterplot of the predictors” included as part of the results of the “Exploratory Data Analysis” reported in <https://onlinecourses.science.psu.edu/stat857/node/224> on the white wine data. These reported results use the full white wine data set $\mathbf{D}_{orig}^{(white)}$, to

construct a matrix of scatterplots of X_i against X_j ; $i \neq j$; $i, j = 1, \dots, 11$. These empirical scatterplots visually suggest stronger correlations between fixed acidity and pH; residual sugar and density; free sulphur dioxide and total sulphur dioxide; density and total sulphur dioxide; density and alcohol—than amongst other pairs of variables. These are the very node pairs that we identify to have edges (at probability in excess of 0.05) between them.

When we compare our learnt graphical model with the results of this reported “Exploratory Data Analysis”, we remind ourselves that partial correlation (that drives the probability of the edge between the i -th and j -th nodes), is often smaller than the correlation between the i -th and j -th variables, computed before the effect of a third variable has been removed (Sheskin, 2004). If this is the case, then an edge between nodes i and j in the learnt graphical model, is indicative of a high correlation between the i -th and j -th variables in the data. However, in the presence of a suppressor variable (that may share a high correlation with the i -th variable, but low correlation with the j -th), the absolute value of the partial correlation parameter can be enhanced to exceed that of the correlation parameter. In such a situation, the edge between the nodes i and j in our learnt graphical model may show up (within our defined 95% HPD credible region on edge probabilities, i.e. at probability higher than 0.05), though the empirical correlation between these variables is computed as low (Sheskin, 2004). So, to summarise, if the empirical correlation between two variables reported for a data set is high, our learnt graphical model should include an edge between the two nodes. But the presence of an edge between pair of nodes is not necessarily an indication of high empirical correlation between a pair of variables—as in cases where suppressor variables are involved. Guessing the effect of such suppressor variables via an examination of the scatterplots is difficult in this multivariate situation. Lastly, it is appreciated that empirical trends are only indicators as to the matrix-Normal density-based model of the learnt correlation structure (and the graphical model learnt thereby) given the data at hand.

Effect on the “quality” variable in the “Exploratory Data Analysis” reported in PSU site, using the white wine data, is examined via a linear regression analysis of the predictors X_1, \dots, X_{11} on the response variable “quality”, which suggests the variables alcohol and volatile acidity to have maximal effect on quality. Indeed, this is corroborated in our learning of the graphical model that manifests edges between the nodes corresponding to variables: alcohol-quality, and volatile acidity-quality.

6.2. Results given data $\mathbf{D}_S^{(red)}$

The $\mathbf{D}_S^{(red)}$ data is the standardised version of a subset of the original red wine data set $\mathbf{D}_{orig}^{(red)}$. $\mathbf{D}_S^{(red)}$ comprises $n = 300$ rows and $p = 12$. The marginal posterior of some of the partial correlation parameters ρ_{ij} computed using the elements of the correlation matrix $\Sigma_S^{(red)}$ (of data $\mathbf{D}_S^{(red)}$) that is updated in

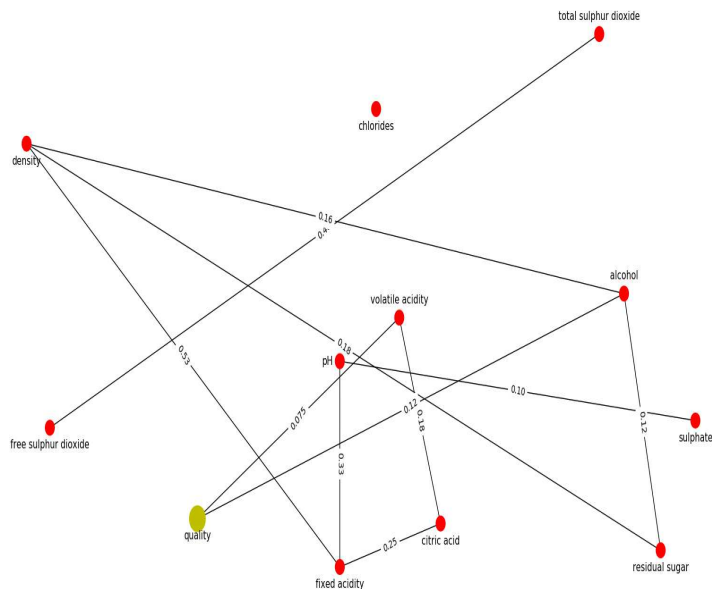


FIG 4. Graphical model of standardised version $\mathbf{D}_S^{(red)}$ of the real data on Portuguese red wine samples (Cortez et al., 1998). Figure is similar to Figure 3, even in the random placement of the nodes, except that this is the graphical model learnt for the red wine data.

the first block of Metropolis-with-2-block-update, are presented in Figure 10 of the Supplementary Material. Figure 11 of the Supplementary Material presents the trace of the joint posterior probability of the edge (G_{ij}) parameters and the variance (v_{ij}) parameters learnt given data $\mathbf{D}_S^{(red)}$. The inferred graphical model of the red wine data is included in Figure 4.

6.2.1. Comparing against empirical work done with red wine data

To the best of our knowledge, analysis of the red wine data has not been reported in the literature. In lieu of that, we undertake construction of a matrix of scatterplots of X_j against X_i from the red wine data. This is shown in Figure 12 of the Supplementary Material, for $i = 1, \dots, 11$. These scatterplots visually indicate moderate correlations between the following pairs of variables: fixed acidity-citric acid, fixed acidity-density, fixed acidity-pH, volatile acidity-citric acid, free sulphur dioxide-total sulphur dioxide, density-alcohol. All these variables share an edge at probability ≥ 0.05 in our learnt graphical model of data $\mathbf{D}_S^{(red)}$ (see Figure 4). We note that all moderately correlated variable pairs, as represented in these scatterplots, are joined by edges in our learnt graphical model of the red wine data – as is to be expected if the learning of the graphical model is correct. Such pairs include fixed acidity-citric acid, fixed acidity-density, fixed acidity-pH, volatile acidity-citric acid, free sulphur dioxide-total sulphur dioxide, density-alcohol. However, an edge may exist between a pair of variables

even when the apparent empirical correlation between these variables is low, owing to the effect of other variables (discussed in Section 6.1.1).

Noticing such edges from the residual-sugar variable, we undertake a regression analysis (ordinary least squares) with residual-sugar regressed against the other remaining 10 vino-chemical variables. The MATLAB output of that analysis carried out using the red wine data, is included in Figure 13 of the Supplementary Material. The analysis indicates that the covariates with maximal (near-equal) effect on residual-sugar, are density and alcohol; indeed, in our learnt graphical model of the red wine data (Figure 4), residual-sugar is noted to enjoy an edge with both density (Z_7) and alcohol (Z_{10})

We also undertook a separate ordinary least squares analysis with the response variable quality, regressed against the vino-chemical variables as the covariates. The MATLAB output of this regression analysis is in Figure 14 of the Supplementary Section. We notice that the strongest (and nearly-equal) effect on quality is from the variables volatile-acidity and alcohol—the very two variables that share an edge at probability ≥ 0.05 with quality, in our learnt graphical model of the red wine data.

7. Metric measuring distance between posterior probability densities of graphs given white and red wine datasets

We seek the distance $\delta(\cdot, \cdot)$ that we defined in Definition 4.4, between the learnt red and white wine graphs, using the method delineated in Section 4. For this, we first compute the normalisation

$S := \max\{\ln(p_{red}^{(0)}), \ln(p_{red}^{(1)}), \dots, \ln(p_{red}^{(N_{iter})}), \ln(p_{white}^{(0)}), \dots, \ln(p_{red}^{(N_{iter})})\}$, which for the red and white wine datasets yields $s = \ln(p_{red}^{(1474)}) \approx 142.7687$. We then use $\exp(\ln(u_m^{(t)})/s)$ in Equation 4.2; $m = white, red$. Then scaling the log posterior given either data set, at any iteration, by the global scale value of $s=142.7687$ approximately, we get $D_H(u_{white}, u_{red}) \approx 0.1153$, so that the logarithm of this value of the Hellinger distance between the 2 learnt graphical models is $\ln(0.1153) \approx -2.1602$. Similarly, using the same scale, the Bhattacharyya distance is $D_B(p_{white}, p_{red}) \approx -1.7623$, where we recall that this measure is a logarithm of the distance.

For this s and the red wine data, we compute the uncertainty inherent in graphical model of the red-wine data as $D_{max,s}(red)$, between the graph that occurs at maximal posterior and that at the minimal posterior (Equation 4.4). Similarly, we compute $D_{max,s}(white)$. We then compute ratio of the Hellinger distance between the graphical models learnt given the red and white-wine data, to the uncertainty inherent in each learnt model, and compare $D_H(p_{white}, p_{red})/D_{max,s}(red)$, with $D_H(p_{white}, p_{red})/D_{max,s}(white)$. This comparison is depicted in the left panel of Figure 5 that shows that the difference $D_{max,s}(white)$ between the scaled posterior of graphs given the white wine data is about 0.0694 while $D_{max,s}(red)$ given the red wine data is about 0.05521, These values are compared to the Hellinger distance (between scaled posteri-

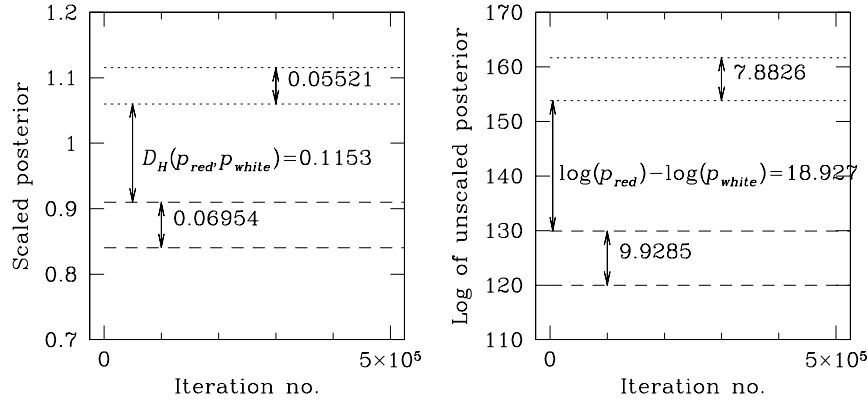


FIG 5. Left: minimum and maximum values of the scaled posterior probability density of the SRGG sampled in an iteration in the MCMC chain run with the red wine data, plotted in dotted lines against the number of the iteration. The difference between these values is depicted within the band delineated by these lines. The broken lines show the same for the results obtained from the MCMC chain run using the white wine data. The value of the Hellinger distance computed using the scaled posterior probabilities of the graphical models given the two wine data sets, is also marked, as about 0.1153. All log posterior values are scaled by a chosen global scale (of about 143), and exponentiated (as discussed in the text). Right: similar to the left panel, except that here, the ratio of the logarithm of the unscaled posteriors is used; the value of the log odds between the posteriors of the red and white wine data sets is marked to be about 18.927.

ors) of about 0.1153, between graphs given the red and white wine data. Thus, $D_H(u_{red}, u_{white})$ is about $1.66D_{max,s}(white)$ and about $2.1D_{max,s}(red)$. Thus, our inter-graph distance metric, between the graphical models learnt given the two data sets is

$$\delta(white, red) \approx 0.44$$

. Then intuitively speaking, this inter-graph distance between the graphical models given the red and white wine datasets, may suggest independence of the data sets.

Again, using the correlation model suggested in Proposition 4.2, the absolute value of the correlation between the 12-dimensional vino-chemical vector-valued measurable for the red wine data and that for the white wine data, is

$$|corr(white, red)| := \exp[-\delta(white, red)] \approx 0.1030,$$

which is a low correlation, indicating that the two graphical models learnt given the real red and white wine Portuguese datasets, are not sampled from the same pdf.

Compared to these, the sample mean of the log odds of the posterior of the graphs generated in the post-burnin iterations, given the two data is 18.9273, which is about 1.9 times the maximal difference between the log posterior values

of graphs achieved in the MCMC run with the white wine data, and about 2.4 times that for the red wine data (see Figure 5). Again, this suggests that the log odds as a measure of divergence between the graphical models given these two wine data sets, is significantly higher than the uncertainty internal to the results for each data.

This clarifies how our pursuit of uncertainties in learnt graphical models, and inter-graph distance, share an integrated umbrage of purpose, where the former leads to the latter.

8. Learning the human disease-symptom network

Our methodology for learning the graphical model, can be implemented even for a highly multivariate data that generates a graph with a very large number of nodes. In this section, we discuss such a graph (with $\gtrsim 8000$ nodes) that describes the correlation structure of the human disease-symptom network.

Hoehndorf et al. (2015) (HSG hereon) learn this network by considering the similarity parameter for each pair of diseases that are elements of an identified set of diseases in the Human Disease Ontology (DO), that contains information about rare and common diseases, and spans heritable, developmental, infectious and environmental diseases. Here, the “similarity parameter” between one disease and another, is computed using the ranked vectors of “normalised pointwise mutual information” (NMPI) parameters for the two diseases, where the NMPI parameter describes the relevance of a symptom (or rather, a phenotype), to the disease in question. HSG define the NMPI parameter semantically, as the normalised number of co-occurrences of a given phenotype and a disease in the titles and abstracts of 5 million articles in Medline. To do this, they make use of the Aber-OWL: Pubmed infrastructure that performs such semantical mining of the Medline abstracts and titles. The disease-disease pairwise semantic similarity parameters – computed using the degree of overlap in the relevance ranks of phenotypes associated with each disease – result in a similarity matrix, which HSG turn into a diseasedisease network based on phenotypes. To do this, they only choose from the top-ranking 0.5% of diseasedisease similarity values. Phenotypes associated with diseases, and corresponding scoring functions (such as the NPMI), exist in the file “doid2hpo-fulltext.txt.gz” at <http://aber-owl.net/aber-owl/diseasephenotypes>. In fact, this file contains information about N_{dis} diseases, and the semantic relevance of each of the N_{pheno} phenotypes to each disease, as quantified by NPMI parameter values, in addition to other scores such as t -scores and z -scores. In this file, N_{dis} is 8676 and N_{pheno} is 19323. In the phenotypic similarity network between diseases that HSG report, diseases are the nodes, and the edge between two nodes exists in this undirected graph, if the similarity between the nodes (diseases) is in the highest-ranking 0.5% of the 38,688,400 similarity values. They remove all self-loops and nodes with a degree of 0. Their network is presented in <http://aber-owl.net/aber-owl/diseasephenotypes/network/>. The network analysis was performed using standard softwares and they identify multi-

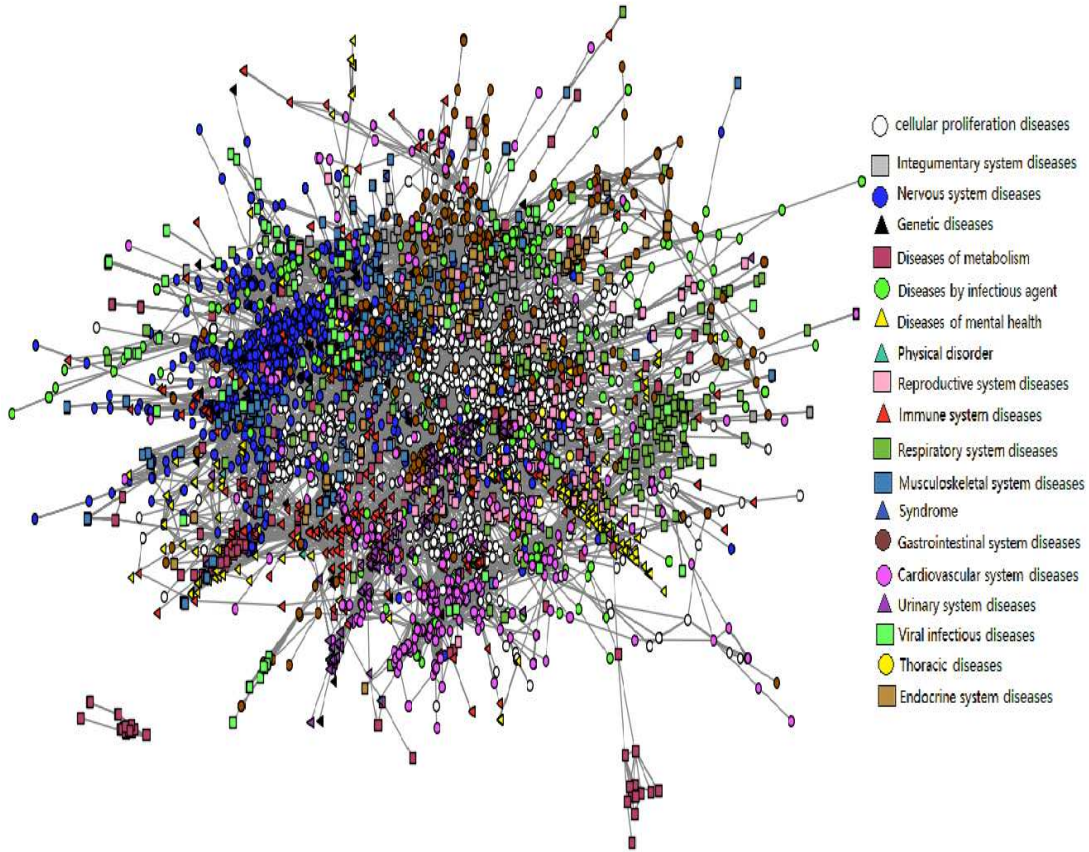


FIG 6. SRGG $G_{m, \mathbf{D}_{DP_h}}(\mathbf{V}, 0.1)$, representing the human disease-phenotype network that we learn using the disease-disease partial correlation obtained using the computed Spearman rank correlation between the rank vectors of a list of phenotypes, where the phenotype ranking reflects semantic relevance of a phenotype to the disease in question (quantified by HSG as the NPMI parameter in the \mathbf{D}_{DP_h} dataset). In our learnt SRGG, $\tau = 0.1$, i.e. only edges (between the i -th and j -th diseases bearing a Spearman rank correlation of $s_{ij}^{(rank)}$), with marginal posterior $m(G_{ij} | \mathcal{S}_{ij}^{(rank)}) \geq 0.9$ are included in this graph. Here cardinality of vertex set \mathbf{V} is 8676, but all nodes with no edges are discarded from this visualised graph, resulting in 6052 diseases (nodes) and 145210 edges remaining that are shown this figure. Diseases identified by HSG, to belong to one of the 19 given disease class, are presented above in the same colour; the colour key identifying these classes, is attached. To draw the graph, we used a Python-based code that implements the Fruchterman-Reingold force-directed algorithm.

ple clusters in their network, with agglomerates of some clusters (of diseases), found to correspond to known disease-classes. The “Group Selector” function on their visualisation kit, allows for the identification of 19 such clusters in their disease-disease network, with each cluster corresponding to a disease-class. To-

tal number of nodes over their identified 19 clusters, is 5059. The number of edges in their network is reported to be 65,795; average node degree ≈ 26.2 .

We discuss detailed comparison of our results to HSG's in the following subsection, including comparison of HSG's and our recovery of the relative number of nodes i.e. diseases, in each of the 19 disease classes that HSG classify their reported network into, and our computed ratios of the averaged intra-class to inter-class variance for each of the 19 classes, compared to the ROC Area Under Curve values reported by HSG for each class.

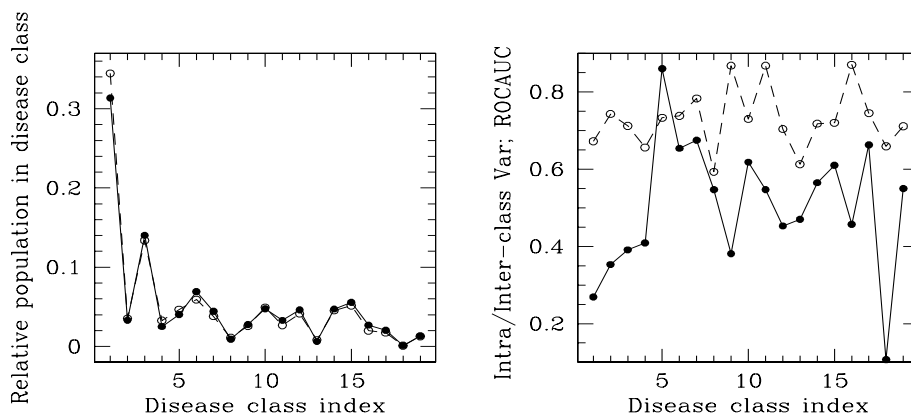


FIG 7. Left: comparison of the relative number of nodes (diseases) that we recover in each of the 19 disease classes that HSG classify their reported network to be classified into, with the relative class-membership reported by HSG. Our results are shown as filled circles joined by solid lines. In open circles threaded by broken lines, we overplot the relative number of diseases in each of the 19 classes, as reported by HSG. Similarity of the relative populations in the different disease classes, indicate that our learnt clustering distribution is similar to that obtained by HSG. Right: our computed ratios of the averaged intra-class to inter-class variance for each of the 19 classes, shown in filled circles; the ROC Area Under Curve values reported by HSG for each class, is overplotted as open circles joined by broken lines. The disease class indices, from assigned values of 1 to 19, are the following respectively: cellular proliferation diseases, integumentary diseases, diseases of the nervous system, genetic diseases, diseases of metabolism, diseases by infectious agents, diseases of mental health, physical disorders, diseases of the reproductive system, of the immune system, of the respiratory system, of the musculoskeletal system, syndromes, gastrointestinal diseases, cardiovascular diseases, urinary diseases, viral infections, thoracic diseases, diseases of the endocrine system.

HSG's network then manifests a similarity-structure that is computed using available NPMI parameter values. Our interest is in learning the disease-disease network as an SRGG, with each edge of such a graphical model learnt to exist at a learnt probability τ . We perform such learning using the NPMI semantic-relevance data that is made available for each of the N_{dis} number of diseases, by HSG; so N_{dis} is the cardinality of the vertex set \mathbf{V} of our sought SRGG. We refer to this human disease-phenotype data as \mathbf{D}_{DP_h} . Using \mathbf{D}_{DP_h} , we first compute the correlation S_{ij} between the i -th and j -th diseases in \mathbf{V} , for each of which, information on the ranked (semantic) relevance of each of the

N_{pheno} phenotypes exist, in this given dataset. Upon computation of pairwise correlations, the SRGG for the data \mathbf{D}_{DP_h} is learnt.

We compute the correlation between the i -th and j -th diseases in the \mathbf{D}_{DP_h} data, $(i, j = 1, \dots, N_{dis}, i \neq j)$, in the following way. We rank the NPMI parameter values for the i -th disease and each of the N_{pheno} phenotypes, with the phenotype of the highest semantic relevance to the i -th disease assigned a rank 1. Let the rank vector of phenotypes, by semantic relevance to the i -th disease take the value τ_i and similarly, the rank vector of phenotypes relevant to the j -th disease is τ_j . We compute the Spearman rank correlation $s_{ij}^{(rank)}$, of vectors τ_i and τ_j . Then we compute this rank correlation $s_{ij}^{(rank)} \forall i, j = 1, \dots, N_{dis}; i \neq j$, between the i -th and j -th nodes. The Spearman rank correlation is preferred to the correlation between the vectors of normalised NPMI values, since we intend to correlate the i -th disease with the j -th disease, depending on how relevant a given list of phenotypes is, to each disease, i.e. depending on the ranked relevance of the phenotypes. We learn the network given this correlation structure, that is itself computed using data \mathbf{D}_{DP_h} (see Section 5 on learning large networks).

Definition 8.1. *Our visualised SRGG in Figure 6 is a sub-graph of the full graph $\mathcal{G}_{m, \mathbf{D}_{DP_h}}(\mathbf{V}, 0.1)$ where \mathbf{V} has cardinality N_{dis} , and the inter-column correlation matrix of data \mathbf{D}_{DP_h} is $\Sigma_C^{(S)} = [s_{ij}^{(rank)}]$, $i \neq j$, $i, j = 1, \dots, N_{dis}$, such that this visualised graph is defined to consist only of nodes with non-zero degree. This visualised graph has 6052 number of nodes (diseases) and 145210 edges, so that the average node degree is about 24. This undirected SRGG represents our learning of the human disease phenotype graph (displayed in Figure 6).*

8.1. Comparing our results to the earlier work done on the human disease-symptom network

The ‘‘Group Selector’’ function on the visualisation kit that HSG use, allows for the identification of 19 such clusters in their disease-disease network, with each cluster corresponding to a disease-class. This function also allows identification of the number of diseases (i.e. nodes) in each disease-class (see left panel of Figure 7). The right panel of Figure 7 displays the ratio of intra-class variance to the inter-class variance of each disease-class; the value of the area under the Receiver Operating Characteristic curve (ROCAUC) for each cluster is overplotted, where the ROCAUC value for the i -th cluster can be interpreted as probability that a randomly chosen node is ranked as more likely to be in the i -th class than in the j -th class; $i \neq j$; $i, j = 1, \dots, 19$ (Hajian-Tilaki, 2013).

9. Conclusion

In this work, we present a methodology for Bayesianly learning a Soft Random Geometric Graph that is drawn in a probabilistic metric space, allowing for the connection function of this SRGG to be equated to the marginal posterior of the graph edge parameter, given the correlation between the points that this edge connects, with the threshold radius on this SRGG to be rendered a probability, s.t. only edges with marginals that exceed such a threshold (probability $\tau \in [0, 1]$) are included in the graph. We demonstrate the SRGG as generated from a point process that we identify as a non-homogenous Poisson process, with intensity that varies with the node.

In fact, correlation between each pair of nodes is learnt as well, and the SRGG updated at each update of the correlation matrix, within each iteration of the iterative inference scheme that we employ; (to be precise, the MCMC-based inference). Here, each of the p nodes of the graph is a variable $X_i - n$ measurements of each of which – comprises the dataset, the standardised version of which we learn the graphical model and the correlation matrix of. To be precise, the i -th column of the dataset contains the n measurements of the r.v. X_i , standardised by its sample mean and standard deviation; $i = 1, \dots, p$. The vertex set of the sought SRGG is then $\mathbf{V} = \{X_1, \dots, X_p\}$. The continuous-valued generative process of the inter-column correlation matrix, is identified after we achieve closed-form marginalisation of the joint likelihood of the inter-column and inter-row correlation matrices given the dataset, over all possible inter-row correlations. The resulting process underlying the inter-column correlation, is then compounded with the non-homogeneous Poisson point process, to generate the SRGG. The graphical model of the data is identified with 95% HPDs, on vertex set \mathbf{V} , to be the graph with edges, the expected marginals of which exceed 5%, where a sample estimate of the expected marginal of an edge is provided by its relative frequency from across the sample of SRGGs that are realised across the iterations of the undertaken inference. When learning a large network, such an iterative inferential scheme is prohibitively expensive though. So then we learn the inter-column correlation of the given dataset empirically, and employ it to learn the SRGG that represents this network.

Our Bayesian learning approach allows for acknowledgement of measurement errors of any observable. The effect of ignoring such existent measurement errors, on the graphical model, is demonstrated using a simple, low-dimensional simulated dataset (see Section 4.2 of the Supplementary Material; compare Figures 4 and 6 of the Supplementary Materials). Even in such a low-dimensional example, the difference made to the inferred graph of the given data, by the inclusion of measurement errors, is clear.

Ultimately we aim at computing the distance between a pair of such learnt graphical models, of respective datasets. To compute this inter-graph distance, we advance a new metric that is given by the difference between the Hellinger distance between the posterior probabilities of the graphs, normalised by the uncertainty in one of the learnt graphs, and the Hellinger distance normalised by the uncertainty in the other learnt graph.

This novel, eventual computation of the inter-graph distance is important in the sense that it informs on the correlation structure of a dataset that is higher-dimensional than being rectangularly-shaped, such as a cuboidally-shaped dataset that comprises slices of rectangularly-shaped data slices. Then, the distance between the graphical models of a pair of such slices of data, will inform on the correlation between such slices of data. Such information is easily calculated under the approach discussed herein, even when the datasets are differently sized, and highly multivariate. An example could be a large network observed on a sample of size n_1 before an intervention/treatment, and after implementation of such intervention, when a smaller sample (of size n_2 ; $n_2 \neq n_1$) is investigated. We illustrate this on computing distance between the learnt vino-chemical graphical models of Portuguese red and white wine samples.

Our learning of large networks is illustrated by the human disease-phenotype network (with $\geq 8,000$ nodes). In this application, learning the inter-node correlation was cast into a semantic exercise in which we learnt the Spearman rank correlation between vectors of associated phenotypes, where any phenotype vector is ranked in order of relevance to the disease in question. Other situations also admit such possibilities, for example, the product-to-product, or service-to-service correlation in terms of associated emotion, (or some other response parameter), can be semantically gleaned from the corpus of customer reviews uploaded to a chosen internet facility, and the same used to learn the network of products/services. Importantly, this method of probabilistic learning of small to large networks, is useful for the construction of networks that evolve with time, i.e. of dynamic networks.

References

- Airoldi, E. M. (2007), “Getting Started in Probabilistic Graphical Models,” *PLoS Computational Biology*, 3(12), e252.
- Bandyopadhyay, D., and Canale, A. (2016), “Sparse Multi-Dimensional Graphical Models: A Unified Bayesian Framework,” *Journal of Royal Statistical Society Series C*, 65(4), 619–640.
- Banerjee, S., Basu, A., Bhattacharya, S., Bose, S., Chakrabarty, D., and Mukherjee, S. (2015), “Minimum distance estimation of Milky Way model parameters and related inference,” *SIAM/ASA Journal on Uncertainty Quantification*, 3(1), 91–115.
- Benner, P., Findeisen, R., Flockerzi, D., Reichl, U., and Sundmacher, K. (2014), *Large-Scale Networks in Engineering and Life Sciences*, Modeling and Simulation in Science, Engineering and Technology, Switzerland: Springer.
- Bhattacharyya, A. (1943), “On a measure of divergence between two statistical populations defined by their probability distributions,” *Bull. Calcutta Math. Soc.*, 35, 99–109.
- Carvalho, C. M., and West, M. (2007), “Dynamic matrix-variate graphical models,” *Bayesian Analysis*, 2(1), 69–97.
URL: <https://doi.org/10.1214/07-BA204>

- Cortez, P., Cerdeira, A., Almeida, F., Matos, T., and Reis, J. (1998), “Modeling wine preferences by data mining from physicochemical properties,” *Decision Support Systems*, 47(4), 547–553.
- Dawid, A. P., and Lauritzen, S. L. (1993), “Hyper-Markov laws in the statistical analysis of decomposable graphical models,” *Ann. Statist.*, 21(3), 1272–1317.
URL: <http://dx.doi.org/10.1214/aos/1176349260>
- Fodor, J. (2004), “Left-continuous t-norms in fuzzy logic: An overview,” *Acta Polytechnica Hungarica* 1(2), 1 (2), 535–537.
- Giles, A. P., Georgiou, O., and Dettmann, C. P. (2016), “Connectivity of Soft Random Geometric Graphs over Annuli,” *Journal of Statistical Physics*, 162(4), 1068–1083.
- Goodman, L. A. (1970), “The Multivariate Analysis of Qualitative Data: Interaction Among Multiple Classifications,” *Journal of the American Statistical Association*, 65, 226–256.
- Gruber, L., and West, M. (2016), “GPU-Accelerated Bayesian Learning and Forecasting in Simultaneous Graphical Dynamic Linear Models,” *Bayesian Analysis*, 11(1), 125–149.
- Guinness, J., Fuentes, M., Hesterberg, D., and Polizzotto, M. (2014), “Multivariate spatial modeling of conditional dependence in microscale soil elemental composition data,” *Spatial Statistics*, 9, 93–108.
- Hajian-Tilaki, K. (2013), “Receiver Operating Characteristic (ROC) Curve Analysis for Medical Diagnostic Test Evaluation,” *Caspian Journal of Internal Medicine*, 4(2), 627–635.
- Hoehndorf, R., Schofield, P. N., and Gkoutos, G. V. (2015), “Analysis of the human diseaseome using phenotype similarity between common, genetic, and infectious diseases,” *Scientific Reports*, 5(10888).
URL: <http://dx.doi.org/10.1038/srep10888>
- Hoff, P. D. (2011), “Separable covariance arrays via the Tucker product, with applications to multivariate relational data,” *Bayesian Analysis*, 6(2), 179–196.
- Kiiveri, H., Speed, T. P., and Carlin, J. B. (1984), “Recursive Causal Models,” *Journal of the Australian Mathematical Society*, 36(Ser.A), 30–52.
- Lauritzen, S. L. (1996), *Graphical Models*, Oxford, UK: Oxford University Press.
- Madigan, D., and Raftery, A. E. (1994), “Model Selection and Accounting for Model Uncertainty in Graphical Models Using Occam’s Window,” *Journal of the American Statistical Association*, 89(428), 1535–1546.
- Matusita, K. (1953), “On the estimation by the minimum distance method,” *Annals of the Institute of Statistical Mathematics*, 5(1), 59–65.
- Menger, K. (1942), “Statistical metrics,” *Proc. Nat. Acad. Sci. USA*, 28 (12), 535–537.
- Ni, Y., Stingo, F. C., and Baladandayuthapani, V. (2017), “Sparse Multi-Dimensional Graphical Models: A Unified Bayesian Framework,” *Journal of the American Statistical Association*, 112(518), 779–793.
- Penrose, M. (2003), *Random Geometric Graphs*, Oxford Studies in Probability, Oxford: OUP.
- Penrose, M. D. (2016), “Connectivity of Soft Tandom Geometric Graphs,” *The*

- Annals of Applied Probability*, 26, 986–1028.
- Robert, C. P., and Casella, G. (2004), *Monte Carlo Statistical Methods*, New York: Springer-Verlag.
- Schweizer, B., and Sklar, A. (1983), *Probabilistic Metric Spaces*, New York: North-Holland.
- Sheskin, D. (2004), *Handbook of parametric and nonparametric statistical procedures*, Boca Raton, Florida: Chapman & Hall/CRC.
- Wang, H., and West, M. (2009), “Bayesian analysis of matrix normal graphical models,” *Biometrika*, 96, 821–834.
- Whittaker, J. (2008), *Graphical Models in Applied Multivariate Statistics*, Switzerland: Wiley.
- Wothke, W. (1993), *Nonpositive definite matrices in structural modeling*, Newbury Park, CA: Sage.
- Xu, Z., Yan, F., and Qi, A. (2012), Infinite tucker decomposition: Nonparametric bayesian models for multiway data analysis,, in *Proceedings of the 29th International Conference on Machine Learning*, pp. 1023–1030.