LEARNING LATENT GRAPH STRUCTURES AND THEIR UNCERTAINTY

Anonymous authors

Paper under double-blind review

ABSTRACT

Graph neural networks use relational information as an inductive bias to enhance prediction performance. Not rarely, task-relevant relations are unknown and graph structure learning approaches have been proposed to learn them from data. Given their latent nature, no graph observations are available to provide a direct training signal to the learnable relations. Therefore, graph topologies are typically learned on the prediction task alongside the other graph neural network parameters. In this paper, we demonstrate that minimizing point-prediction losses does not guarantee proper learning of the latent relational information and its associated uncertainty. Conversely, we prove that suitable loss functions on the stochastic model outputs simultaneously grant solving two tasks: (i) learning the unknown distribution of the latent graph and (ii) achieving optimal predictions of the model output. Finally, we propose a sampling-based method that solves this joint learning task. Empirical results validate our theoretical claims and demonstrate the effectiveness of the proposed approach.

024 025 026

027

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023

1 INTRODUCTION

028 Relational information processing has provided breakthroughs in the analysis of rich and complex 029 data coming from, e.g., social networks, natural language, and biology. This side information takes various forms, from structuring the data into clusters to defining causal relations and hierarchies, and enables machine learning models to condition their predictions on dependency-related observations. 031 In this context, predictive models take the form $y = f_{\psi}(x, A)$, where the input-output relation $x\mapsto y$ – modeled by f_ψ and its parameters in ψ – is conditioned on the relational information 033 encoded in variable A. Graph Neural Networks (GNNs) [Scarselli et al., 2008] are one example 034 of models of this kind that rely on a graph structure represented as an adjacency matrix A and have been demonstrated successful in a plethora of applications [Fout et al., 2017; Shlomi et al., 2020]. Throughout this paper, we focus on predictors where A is an adjacency matrix, although the 037 theoretical results we develop are valid for A being any discrete latent random variable.

Indeed, relational information is needed to implement such a relational inductive bias and, in some cases, it is provided at the application design phase. However, more frequently, such topological information is not rich enough to address the problem at hand, and – not seldom – it is completely unavailable. Therefore, Graph Structure Learning (GSL) emerges as an approach to learn the graph topology [Kipf et al., 2018; Franceschi et al., 2019; Yu et al., 2021; Fatemi et al., 2021; Zhu et al., 2021; Cini et al., 2023] alongside the predictive model f_{ψ} . This entails formulating a joint learning process that learns the adjacency matrix A – or a parameterization of it – along with the predictor's parameters ψ . This can be achieved by optimizing a loss function, e.g., a point prediction measure based on the square or the absolute prediction error.

Different sources of uncertainty affect the graph structure learning process, including epistemic uncertainty in the data and variability inherent in the data-generating process. Learning appropriate models of the data-generating process can provide valuable insights into the modeled environment with uncertainty quantification enhancing explainability and interpretability, ultimately enabling more informed decision-making. Examples of applications are found in the study of infection and information spreading, as well as biological systems [Gomez Rodriguez et al., 2013; Lokhov, 2016; Deleu et al., 2022]. It follows that a probabilistic framework is appropriate to accurately capture the uncertainty in the learned relations whenever randomness affects the graph topology. Probabilistic

approaches have been devised in recent years. For instance, research carried out in Franceschi et al. [2019]; Zhang et al. [2019]; Elinas et al. [2020]; Cini et al. [2023] propose methods that learn a parametric distribution P_A^{θ} over the latent graph structure A. However, none of them have studied whether these approaches were able to learn a *calibrated* latent distribution P_A^{θ} , properly reflecting the uncertainty associated with the learned topology.

In this paper, we address the joint problem of learning a predictive model yielding optimal pointprediction performance of the output y and, contextually, a calibrated distribution for the latent adjacency matrix A. In particular, the novel contributions can be summarized as:

- 1. We demonstrate that models trained to achieve optimal point predictions do *not* guarantee calibration of the adjacency matrix distribution [Section 4].
- 2. We provide theoretical conditions on the predictive model and loss function that guarantee both distribution calibration and optimal point-predictions [Section 5].
- 3. We propose a theoretically grounded sampling-based learning method to address the joint learning problem [Section 5].
- 4. We empirically validate the theoretical developments and claims presented in this paper and show that the proposed method is indeed able to solve the joint learning task [Section 6].

2 RELATED WORK

063

064

065

066

067

068

069

071

073

074

075 **Graph Structure Learning** GSL is often employed end-to-end with a predictive model to better 076 solve a downstream task. Examples include applications within graph deep learning methods for 077 static [Jiang et al., 2019; Yu et al., 2021; Kazi et al., 2022] and temporal data [Wu et al., 2019; 2020; Cini et al., 2023; De Felice et al., 2024]; a recent review is provided by Zhu et al. [2021]. Some approaches from the literature model the latent graph structure as stochastic [Kipf et al., 2018; 079 Franceschi et al., 2019; Elinas et al., 2020; Shang et al., 2021; Cini et al., 2023], mainly as a way to enforce sparsity of the adjacency matrix. To operate on discrete latent random variables, Franceschi 081 et al. [2019] utilize straight-through gradient estimations, Cini et al. [2023] rely on score-based 082 gradient estimators, while Niepert et al. [2021] design an implicit maximum likelihood estimation 083 strategy. A different line of research is rooted in graph signal processing, where the graph is es-084 timated from a constrained optimization problem and the smoothness assumption of the signals 085 [Kalofolias, 2016; Dong et al., 2016; Mateos et al., 2019; Coutino et al., 2020; Pu et al., 2021]. A few works from the Bayesian literature have tackled the task of estimating uncertainties associated 087 with graph edges. The model-based approaches by Lokhov [2016]; Gray et al. [2020] are two ex-880 amples tackling relevant applications benefiting from uncertainty quantification. Within the deep learning literature, Zhang et al. [2019] propose a Bayesian Neural Network (BNN) modeling the 089 random graph realizations. Differently, Wasserman & Mateos [2024] develop a BNN designed over 090 graph signal processing principles. While some results on the output calibration are exhibited, to 091 the best of our knowledge, no guarantee or evidence of calibration of the latent variable is provided, 092 which we study in this paper instead.

094 **Calibration of the model's output** Research on model calibration has primarily focused on obtaining accurate and consistent predictions of the statistical properties of the target (random) vari-096 ables y, from which uncertainty estimates on the model's predictions are derived. For discrete outputs, such as in classification tasks, Guo et al. [2017] investigated the calibration of modern deep 097 learning models and proposed temperature scaling as a solution. Other techniques in the same con-098 text include Histogram Binning [Zadrozny & Elkan, 2001], Cross Entropy loss with label smoothing [Müller et al., 2019], and Focal Loss [Mukhoti et al., 2020]. For continuous output distributions, 100 Laves et al. [2020] proposed σ scaling, while Kuleshov et al. [2018] developed a technique inspired 101 by Platt scaling. More recently, conformal prediction techniques [Shafer & Vovk, 2008] have gained 102 popularity for providing confidence intervals in predictions. We stress that within this paper, we are 103 mainly concerned with latent variable calibration, rather than output calibration, although the two 104 are related to each other. 105

Deep latent variable models Latent variables are extensively used in deep generative modeling
 [Kingma & Welling, 2013; Rezende et al., 2014], both with continuous and discrete latent variables
 [Van Den Oord et al., 2017; Bartler et al., 2019]. In deep models, latent random variables often

108 lack direct physical meaning, with only the outputs being collected for training. Therefore, studies 109 mainly focused on maximizing the likelihood of the observed outputs in the training set, rather 110 than calibrating the latent distribution. A few works proposed regularization of the latent space to 111 improve stability and accuracy [Xu & Durrett, 2018; Joo et al., 2020], facilitate smoother transitions 112 in the output when the latent variable is slightly modified [Hadjeres et al., 2017], and apply other techniques aimed at enhancing data generation or improving model performance in general [Connor 113 et al., 2021]. 114

115 To the best of our knowledge, no prior work has studied the joint learning problem of calibrating the 116 latent graph distribution while achieving optimal point predictions.

117 118 119

120

125 126

127

128 129

134 135

153

154

157 158

3 PROBLEM FORMULATION

Consider a set of N interacting entities and the data-generating process

$$\begin{cases} A \sim P_A^* \\ y = f^*(x, A) \end{cases}$$
(1)

where $y \in \mathcal{Y}$ is the system output obtained from input $x \in \mathcal{X}$ through function f^* and conditioned on a realization of the latent adjacency matrix $A \in \mathcal{A} \subseteq \{0,1\}^{N \times N}$ drawn from distribution P_A^* ; input x is assumed to be drawn from any distribution P_x^* and superscript * refers to unknown entities. Each entry of the adjacency matrix A is a binary value encoding the existence of a pairwise relation between two nodes. In the sequel, $x \in \mathcal{X} \subseteq \mathbb{R}^{N \times d_{in}}$ and $y \in \mathcal{Y} \subseteq \mathbb{R}^{N \times d_{out}}$ are stacks of N node-level feature vectors of dimension d_{in} and d_{out} , respectively, representing continuous inputs 130 and outputs.

131 Given a training dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ of n input-output observations from (1), we aim at 132 learning a probabilistic predictive model 133

$$\begin{cases} A \sim P_A^{\theta} \\ \hat{y} = f_{\psi}(x, A) \end{cases}$$
(2)

136 from \mathcal{D} , while learning at the same time distribution P_A^{θ} approximating P_A^* . The two parameter 137 vectors θ and ψ are trained to approximate distinct entities in (1), namely the distribution P_A^* and 138 function f^* , respectively. We assume 139

Assumption 3.1. The family $\{P_A^\theta\}$ of probability distributions P_A^θ parametrized by θ and the family 140 of predictive functions $\{f_{\psi}\}$ are expressive enough to contain the true latent distribution P_{A}^{*} and 141 function f^* , respectively. 142

143 Assumption 3.1 implies that $f^* \in \{f_{\psi}\}$ and $P^*_A \in \{P^{\theta}_A\}$ but does not request uniqueness of the parameters vectors ψ^* and θ^* such that $f_{\psi^*} = f^*$ and $P^{\theta^*}_A = P^*_A$. Under such assumption the mini-144 145 mum function approximation error is null and we can focus on the theoretical conditions requested 146 to guarantee successful learning, i.e., achieving both optimal point predictions and latent distribution 147 calibration. In Section 6.2, we empirically show that the theoretical results can extend beyond this assumption in practice. 148

149 **Optimal point predictions** Outputs y and \hat{y} of probabilistic model (1) and (2) are random vari-150 ables following push-forward distributions¹ $P_{y|x}^*$ and $P_{y|x}^{\theta,\psi}$, respectively. A single point prediction 151 $y_{PP} \in \mathcal{Y}$ can be obtained through an appropriate functional $T[\cdot]$ as 152

$$y_{PP} = y_{PP}(x,\theta,\psi) \equiv T \left[P_{y|x}^{\theta,\psi} \right].$$
(3)

155 For example, T can be the expected value or the value at a specific quantile. We then define an optimal predictor as one whose parameters θ and ψ minimize the expected point-prediction loss 156

$$\mathcal{L}^{point}(\theta,\psi) = \mathbb{E}_{x \sim P_x^*} \left[\mathbb{E}_{y \sim P_{y|x}^*} \left[\ell \left(y, y_{PP}(x,\theta,\psi) \right) \right] \right]$$
(4)

159 between the system output y and the point-prediction y_{PP} , as measured by of a loss function ℓ : 160 $\mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_+.$ 161

¹The distribution of $y = f^*(x, A)$ originated from P^*_A and of $\hat{y} = f_{\psi}(x, A)$ originated from P^{θ}_A .

Statistical functional T is coupled with the loss ℓ as the optimal functional T to employ given a specific loss ℓ is often known [Berger, 1990; Gneiting, 2011], when $P_{y|x}^{\theta,\psi}$ approximates well $P_{y|x}^*$. For instance, if ℓ is the Mean Absolute Error (MAE) the associated functional T is the median, if ℓ is the Mean Squared Error (MSE) the associated functional is the expected value.

167 Latent distribution calibration Calibration of a parametrized distribution P_A^{θ} requires learning 168 parameters θ , so that P_A^{θ} aligns with true distribution P_A^* . Quantitatively, a dissimilarity measure 169 $\Delta^{cal} : \mathcal{P}_A \times \mathcal{P}_A \to \mathbb{R}_+$, defined over a set \mathcal{P}_A of distributions on \mathcal{A} , assesses how close two 170 distributions are. The family of *f*-divergences [Rényi, 1961], such as the Kullback-Leibler diver-171 gence, and the integral probability metrics [Müller, 1997], such as the maximum mean discrepancy 172 [Gretton et al., 2012] are examples of such dissimilarity measures. In this paper, we are interested 173 in those discrepancies for which $\Delta^{cal}(P_1, P_2) = 0 \iff P_1 = P_2$ holds. It follows that the latent 174 distribution P_A^{θ} is calibrated on P_A^* if it minimizes the latent distribution loss

175

183

184 185

186

187

188

189 190

191 192 $\mathcal{L}^{cal} = \mathbb{E}_{x \sim P_x^*} \left[\Delta^{cal} \left(P_A^*, P_A^\theta \right) \right], \tag{5}$

or simply $\mathcal{L}^{cal} = \Delta^{cal} \left(P_A^*, P_A^{\theta} \right)$, when A and x are independent.

The problem of designing a predictive model (2) that both yields optimal point predictions (i.e., minimizes \mathcal{L}^{point} in (4)) and calibrates the latent distribution (i.e., minimizes \mathcal{L}^{cal} in (5)) is nontrivial for two main reasons. At first, as the latent distribution P_A^* is unknown (and no samples from it are available), we cannot directly estimate \mathcal{L}^{cal} . Second, as shown in Section 4, multiple sets of θ parameters may minimize \mathcal{L}^{point} without minimizing \mathcal{L}^{cal} .

4 LIMITATIONS OF POINT-PREDICTION OPTIMIZATION

In this section, we demonstrate that the optimization of a point prediction loss (Equation (4)) does not generally grant calibration of the latent random variable A.

Proposition 4.1. Consider Assumption 3.1. Loss function $\mathcal{L}^{point}(\theta, \psi)$ in (4) is minimized by all θ and ψ s.t. $T\left[P_{y|x}^{\theta,\psi}\right] = T\left[P_{y|x}^*\right]$ almost surely on x and, in particular,

$$P_{y|x}^{point}(\theta,\psi) \text{ is minimal } \stackrel{\longleftarrow}{\Longrightarrow} P_{y|x}^{\theta,\psi} = P_{y|x}^*.$$

The proof of the proposition is given in Appendix A.1; we provide a counterexample for which calibration is not granted even when the processing function f_{ψ} is equal to f^* in Appendix A.2.

L

Figure 1 empirically demonstrates that optimizing point prediction losses does not necessarily 199 guarantee distribution calibration. In particular, 200 we compute different losses between data gen-201 erated with a ground truth system model (model 202 (1) with optimal parameter θ^*) and outputs pro-203 duced with a different model (model (2), with 204 varying θ parameters). In red, the MAE is used 205 as the loss function ℓ in the point prediction loss 206 \mathcal{L}^{point} of (4). Since all $\theta \geq 0.725$ produce statistically equivalent losses, this simple experi-207 ment demonstrates the inefficacy of minimiz-208 ing \mathcal{L}^{point} for latent distribution calibration. In 209 blue, we show the loss we propose in the next 210 section, which presents a minimum in θ^* . The 211



Figure 1: A data generating model, as in (1), is used to produce a dataset with latent distribution parameter θ^* . Outputs are generated for different values θ as in (2). In red, losses are computed as in (4) with ℓ being the MAE. In blue, losses are computed with our approach described further on.

details of this experiment can be found in Section 6.1. However, we recommend reading the entire paper first to better understand the experiment's context and setting.

Given the provided negative result and the impossibility of assessing loss \mathcal{L}^{cal} in (5), in the next section, we propose another optimization objective that, as we will prove, allows us to both calibrate the latent random variable and to have optimal point predictions.

5 PREDICTIVE DISTRIBUTION OPTIMIZATION: TWO BIRDS WITH ONE STONE

218 In this section, we show that we can achieve an optimal point predictor (2) and a calibrated latent distribution P_A^{θ} by comparing push-forward distributions $P_{y|x}^*$ and $P_{y|x}^{\theta,\psi}$ of the outputs y conditioned 219 220 on input x. In particular, Theorem 5.2 below proves that, under appropriate conditions, minimization 221 of the output distribution loss 222

223 224

225 226

227

229

230

234 235

216

217

 $\mathcal{L}^{dist}(\theta,\psi) = \mathbb{E}_{x \sim P_x^*} \left[\Delta(P_{y|x}^*, P_{y|x}^{\theta,\psi}) \right]$ (6)

provides calibrated P_A^{θ} , even when P_A^* is not available; $\Delta : \mathcal{P}_y \times \mathcal{P}_y \to \mathbb{R}_+$ is a dissimilarity measure between distributions over space \mathcal{Y} . We assume the following on dissimilarity measure Δ . Assumption 5.1. $\Delta(P_1, P_2) \ge 0$ for all distributions P_1 and P_2 in \mathcal{P}_y and $\Delta(P_1, P_2) = 0$ if and only if $P_1 = P_2$. 228

Several choices of Δ meet Assumption 5.1, e.g., f-divergences and some integral probability metrics [Müller, 1997]; the dissimilarity measure Δ employed in this paper is discussed in Section 5.1.

231 **Theorem 5.2.** Let $I = \{x : A \mapsto f^*(x, A) \text{ is injective}\} \subseteq \mathcal{X}$ be the set of points $x \in \mathcal{X}$ such that 232 map $A \mapsto f^*(x, A)$ is injective. Under Assumptions 3.1 and 5.1, if $\mathbb{P}_{x \sim P_x^*}(I) > 0$, then 233

$$\mathcal{L}^{dist}(\theta,\psi^*) = 0 \implies \begin{cases} \mathcal{L}^{point}(\theta,\psi^*) \text{ is minimal} \\ \mathcal{L}^{cal}(\theta) = 0, \end{cases}$$

236 where ψ^* is such that $f_{\psi^*} = f^*$. 237

238 Theorem 5.2 is proven in Appendix A.3. Under the theorem's hypotheses, a predictor that mini-239 mizes \mathcal{L}^{dist} is both *calibrated* on the latent random distribution and provides *optimal point predictions*. This overcomes limits of Proposition 4.1 where optimization of $\mathcal{L}^{point}(\theta, \psi^*)$ does not grant 240 $\mathcal{L}^{cal}(\theta) = 0.$ 241

242 The hypotheses under which Theorem 5.2 holds are rather mild. In fact, condition $\mathbb{P}_{r \sim P^*}(I) > 0$ 243 pertains to the data-generating process and intuitively ensures that, for some x, different latent ran-244 dom variables produce different outputs. A sufficient condition for $\mathbb{P}_{x \sim P_x}(I) > 0$ to hold is the 245 existence of a point \bar{x} in the support of P_x^* such that $A \mapsto f^*(\bar{x}, A)$ is injective with f^* continuous w.r.t. \bar{x} ; see Corollary A.1 in Appendix A.3. Although only a single point \bar{x} is required, having more 246 points that satisfy the condition simplifies the training of the parameters. Corollary A.1 holds for 247 arbitrarily complex processing functions f^* . More specifically, when considering simple GNN lay-248 ers and discrete latent matrices A, we can prove that the condition $\mathbb{P}_{x \sim P^*}(I) > 0$ is – except from 249 pathological cases – always satisfied (see Proposition A.2 in Appendix A.3). Instead, condition 250 $f_{\psi} = f^*$ is set to avoid scenarios of different, yet equivalent,² representations of the latent distribu-251 tion. An empirical analysis of the theorem's assumptions is provided in Section 6.2, demonstrating that the theoretical results hold in practice, even when the assumption does not strictly apply. 253

Assumptions 3.1 and 5.1 can be met with an appropriate choice of model (2) and measure Δ ; as 254 such they are controllable by the designer. Assumption 5.1 prevents from obtaining mismatched 255 output distributions when $\mathcal{L}^{dist}(\theta, \psi) = 0$ and can be easily satisfied. As mentioned above, popular 256 measures, e.g., the Kullback-Leibler divergence, meet the theorem's assumptions and therefore can 257 be adopted as Δ . However, as f-divergences rely on the explicit evaluation of the likelihood of y, 258 they are not always practical to compute [Mohamed & Lakshminarayanan, 2016]. For this reason, 259 we propose considering the Maximum Mean Discrepancy (MMD) [Gretton et al., 2012] as a versa-260 tile alternative that allows Monte Carlo computation without requiring evaluations of the likelihood w.r.t. the output distributions $P_{y|x}^*$ and $P_{y|x}^{\theta,\psi}$. Energy distances [Székely & Rizzo, 2013] provide an 261 262 alternative feasible choice.

263 264 265

266

267

268 269

5.1 MAXIMUM MEAN DISCREPANCY

Given two distributions $P_1, P_2 \in \mathcal{P}_y$, MMD can be defined as

$$\mathbf{MMD}_{\mathcal{G}}[P_1, P_2] = \sup_{g \in \mathcal{G}} \left\{ \mathbb{E}_{y \sim P_1} \left[g(y) \right] - \mathbb{E}_{y \sim P_2} \left[g(y) \right] \right\},\tag{7}$$

²E.g., $f_{\psi}(A, x) = f_*(\mathbf{1} - A, x)$ and P_A^{θ} encoding the absence of edges instead of their presence as in P_A^* .

273 274 275

276

277

278

279

281 282

283 284

285 286

296

314

i.e., the supremum, taken over a set \mathcal{G} of functions $\mathcal{Y} \to \mathbb{R}$, of the difference between expected values w.r.t. P_1 and P_2 . An equivalent form is derived for a generic kernel function $\kappa(\cdot, \cdot) : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$:

$$\mathbf{MMD}_{\mathcal{G}_{\kappa}}^{2}[P_{1}, P_{2}] = \mathbb{E}_{\substack{y_{1}, y_{1}^{\prime} \sim P_{1}}} \left[\kappa(y_{1}, y_{1}^{\prime}) \right] - 2\mathbb{E}_{\substack{y_{1} \sim P_{1} \\ y_{2} \sim P_{2}}} \left[\kappa(y_{1}, y_{2}) \right] + \mathbb{E}_{\substack{y_{2}, y_{2}^{\prime} \sim P_{2}}} \left[\kappa(y_{2}, y_{2}^{\prime}) \right]$$
(8)

and it is associated with the unit-ball \mathcal{G}_k of functions in the reproducing kernel Hilbert space of κ ; note that (8) is the square of (7). Moreover, when universal kernels are considered (e.g., the Gaussian one), then (8) fulfills Assumption 5.1 (see Theorem 5 of Gretton et al. [2012]). Dissimilarity in (8) can be conveniently estimated via Monte Carlo (MC) and employed within a deep learning framework. Accordingly, we set $\Delta = \text{MMD}_{\mathcal{G}_{\kappa}}^2$ and learn parameter vectors ψ and θ by minimizing $\mathcal{L}^{dist}(\theta, \psi)$ via gradient-descent methods.

5.2 FINITE-SAMPLE COMPUTATION OF THE LOSS

To compute the gradient of $\mathcal{L}^{dist}(\theta, \psi) = \mathbb{E}_{x \sim P_x^*} \left[\text{MMD}_{\mathcal{G}_{\kappa}}^2 \left[P_{y|x}^{\theta,\psi}, P_{y|x}^* \right] \right]$ w.r.t. parameter vectors ψ and θ , we rely on MC sampling to estimate in (6) expectations over input $x \sim P_x^*$, target output $y \sim P_{y|x}^{\theta,\psi}$ and model output $\hat{y} \sim P_{y|x}^{\theta,\psi}$. This amounts to substitute $\text{MMD}_{\mathcal{G}_{\kappa}}^2$ with

$$\widehat{\text{MMD}}^{2}(\theta,\psi;x,y) = \frac{\sum_{i=1}^{N_{adj}} \sum_{j=1}^{i-1} \kappa(\hat{y}_{i},\hat{y}_{j})}{N_{adj}(N_{adj}-1)} - 2\frac{\sum_{i=1}^{N_{adj}} \kappa(y,\hat{y}_{i})}{N_{adj}}$$
(9)

In (9), $N_{adj} > 1$ is the number of adjacency matrices sampled from P_A^{θ} to obtain output samples $\hat{y}_i = f_{\psi}(x, A_i) \sim P_{y|x}^{\theta, \psi}$, whereas the pair (x, y) is a pair from the training set \mathcal{D} . We remark that in (9) the third term of (8) – i.e., the one associated with the double expectation with respect to $P_{y|x}^*$ – is neglected as it does not depend on ψ and θ .

Gradient $\nabla_{\psi} \mathcal{L}^{dist}(\theta, \psi)$ is computed via automatic differentiation by averaging $\nabla_{\psi} \widehat{\text{MMD}}^2(\theta, \psi)$ within a mini-batch of observed data pairs $(x_i, y_i) \in \mathcal{D}$. For $\nabla_{\theta} \mathcal{L}^{dist}(\theta, \psi)$, the same approach is not feasible. This limitation arises because the gradient is computed with respect to the same parameter vector θ that defines the integrated distribution. Here, we rely on a score-function gradient estimator (SFE) [Williams, 1992; Mohamed et al., 2020] which uses the log derivative trick to rewrite the gradient of an expected loss L as $\nabla_{\theta} \mathbb{E}_{A \sim P^{\theta}}[L(A)] = \mathbb{E}_{A \sim P^{\theta}}[L(A)\nabla_{\theta} \log P^{\theta}(A)]$, with $P^{\theta}(A)$ denoting the likelihood of $A \sim P^{\theta}$. Applying the SFE to our problem the gradient of the loss function w.r.t. θ reads:

$$\nabla_{\theta} \mathcal{L}^{dist}(\psi, \theta) = \mathbb{E}_{(x, y^*) \sim P_{x, y}^*} \left[\mathbb{E}_{\hat{y}_1, \hat{y}_2 \sim P_{y|x}^{\theta, \psi}} \left[\kappa(\hat{y}_1, \hat{y}_2) \nabla_{\theta} \log \left(P_{y|x}^{\theta, \psi}(\hat{y}_1) P_{y|x}^{\theta, \psi}(\hat{y}_2) \right) \right] - 2 \mathbb{E}_{\hat{y} \sim P_{y|x}^{\theta, \psi}} \left[\kappa(y^*, \hat{y}) \nabla_{\theta} \log P_{y|x}^{\theta, \psi}(\hat{y}) \right] \right]$$
(10)

An apparent setback of SFEs is their high variance [Mohamed et al., 2020], which we address in Section 5.3 by deriving a variance-reduction technique based on control variates that requires negligible computational overhead.

315 5.3 VARIANCE-REDUCED LOSS FOR SFE 316

Two natural approaches to reduce the variance of MC estimates of (10) involve (i) increasing the number B of training data points in the mini-batch used for each gradient estimate and (ii) increasing the number N_{adj} of adjacency matrices sampled for each data point in (9). These techniques act on two different sources of noise. Increasing B decreases the variance coming from the data-generating process, whereas increasing N_{adj} improves the approximation of the predictive distribution $P_{y|x}^{\theta,\psi}$. Nonetheless, by fixing B and N_{adj} , it is possible to further reduce the latter source of variance by employing the *control variates* method [Mohamed et al., 2020] that, in our case, requires only a negligible computational overhead but sensibly improves the training speed (see Section 6). Consider the expectation $\mathbb{E}_{A \sim P^{\theta}}[L(A) \nabla_{\theta} \log P^{\theta}(A)]$ of the SFE – both terms in (10) can be cast into that form. With the control variates method, L(A) is replaced by a surrogate function

$$\tilde{L}(A) = L(A) - \beta \Big(h(A) - \mathbb{E}_{A \sim P^{\theta}}[h(A)] \Big)$$
(11)

that leads to a reduced variance in the MC estimator while maintaining it unbiased. In this paper, we set function h(A) to $\nabla_{\theta} \log P^{\theta}(A)$ and show how to compute a near-optimal choice for scalar value β , often called *baseline* in the literature. As the expected value of $\nabla_{\theta} \log P^{\theta}(A)$ is zero, gradient (10) rewrites as

$$\nabla_{\theta} \mathcal{L}^{dist} = \mathbb{E}_{(x,y^*) \sim P_{x,y}^*} \left[\mathbb{E}_{A_1, A_2 \sim P_A^{\theta}} \left[(\kappa(f_{\psi}(x, A_1), f_{\psi}(x, A_2)) - \beta_1) \ \nabla_{\theta} \log\left(P_A^{\theta}(A_1)P_A^{\theta}(A_2)\right) \right] - 2 \mathbb{E}_{A \sim P_A^{\theta}} \left[(\kappa(y^*, f_{\psi}(x, A)) - \beta_2) \ \nabla_{\theta} \log P_A^{\theta}(A) \right] \right].$$
(12)

336 337 338

343 344

345

359

360

327 328

330

331

In Appendix **B**, we show that in our setup the best values of β_1 and β_2 are approximated by

$$\tilde{\beta}_1 = \mathbb{E}_{\substack{x \sim P_x^*\\A_1, A_2 \sim P_A^\theta}} \left[\kappa \left(f_\psi(x, A_1), f_\psi(x, A_2) \right) \right], \qquad \tilde{\beta}_2 = \mathbb{E}_{\substack{(x, y^*) \sim P_{x, y}^*\\A \sim P_A^\theta}} \left[\kappa \left(y^*, f_\psi(x, A) \right) \right], \tag{13}$$

which can be efficiently computed via MC reusing the kernel values already computed for (12).

5.4 COMPUTATIONAL COMPLEXITY

346 Focusing on the most significant terms, for every data pair (x, y) in the training set, computing the 347 loss \mathcal{L}^{dist} requires $\mathcal{O}(N_{adj}^2)$ kernel evaluations $\kappa(\hat{y}_i, \hat{y}_j)$ in (9), $\mathcal{O}(N_{adj})$ forward passes through the 348 GNN $\hat{y}_i = f_{\psi}(x, A_i)$ in (9) and $\mathcal{O}(N_{adj})$ likelihood computations $P_A^{\theta}(A_i)$ in (12). The computation 349 of baselines β_1 and β_2 in (13) requires virtually no overhead, as commented in previous Section 5.3. Similarly, computing the loss gradients requires $\mathcal{O}(N_{adj}^2)$ derivatives for what concerns the kernels, 350 $\mathcal{O}(N_{adj})$ gradients $\nabla_{\psi}\hat{y}_i$ and $\nabla_{\theta}\log P^{\theta}_A(A_i)$. We empirically observed that for $N_{adj} \geq 16$, both the 351 352 latent distribution loss \mathcal{L}^{cal} and the point prediction loss \mathcal{L}^{point} of final models are equivalent for 353 the considered problem. This suggests that N_{adj} is not a critical hyperparameter.

Since we can employ sparse representations of adjacency matrices, the GNN processing costs scale linearly in the number of nodes N for bounded-degree graphs. From our experience, the GNN processing is the most demanding operation and the cost of quadratic components, such as the parameterization of θ_{ij} , do not pose significant overhead.

6 EXPERIMENTS

361 This section empirically validates the proposed technique and the main claims of the paper. To assess 362 the calibration performance of models, it is necessary to compare the learned graph distribution P_A^{θ} with the ground-truth latent distribution P_A^* . However, to our knowledge, no real-world datasets 364 provide such ground truth. Therefore, we developed theoretical guarantees to support the application 365 of these methods to real data and - in this section - we conduct the empirical validation using 366 synthetic data. Section 6.1 demonstrates that the proposed approach can successfully solve the joint 367 learning problem across different graph sizes, highlights the effectiveness of the variance reduction 368 technique, and reveals challenges in optimizing point prediction losses when also aiming for latent variable calibration. Section 6.2 empirically investigates the generality of the theoretical results we 369 develop, demonstrating appropriate calibration of the latent distribution even in scenarios where the 370 assumptions of Theorem 5.2 do not hold. 371

In order to assess the latent variable calibration performance, i.e., the discrepancy between P_A^* and the learned P_A^{θ} , the ground-truth latent distribution P_A^* must be given. Such ground-truth knowledge is not available in real-world applications as the latent distribution is indeed unknown. For this reason, we designed a synthetic dataset that allows us to evaluate different performance metrics on both y and A while controlling several properties of the task, like the number of nodes and the probability of each edge. We remark that the latent distribution P_A^* is used *only* to assess performance and does not drive the model training in any way. 389

390

391 392

416

428



Figure 2: Validation losses \mathcal{L}^{dist} , \mathcal{L}^{cal} and \mathcal{L}^{point} during training. At epoch 5, the learning rate is decreased to ensure convergence. \mathcal{L}^{dist} in Subfigure 2a is negative as the third term in (8) is constant and not considered.

393 **Dataset and models** Consider data-generating process (1) with latent distribution $P_A^* = P_A^{\theta^*}$ 394 producing N-node adjacency matrices. P_A^* is defined by a set of $N \times N$ independent Bernoulli 395 distributions, each of which corresponds to the sampling probability of an edge. Function $f_* = f_{\psi^*}$ is a generic GNN with node-level readout, i.e., $f_{\psi^*}(\cdot, A) : \mathbb{R}^{N \times d_{in}} \to \mathbb{R}^{N \times d_{out}}$. In the below 397 experiments, N is set to 12, while input and output node feature dimensions are $d_{in} = 4$ and 398 $d_{out} = 1$, respectively. The components θ^* are set to either 0 or 3/4 according to the pattern 399 depicted in Figure 7; the specifics of f_{ψ^*} and P_x^* are detailed in Appendix C. We result in a dataset 400 of 35k input-output pairs (x, y), 80% of which are used as training set, 10% as validation set, and 401 the remaining 10% as test set. As predictive model family (2), we follow the same architecture 402 of f_{ψ^*} and $P_A^{\theta^*}$ ensuring that during all the experiments Assumption 3.1 is fulfilled. The model 403 parameters are trained by optimizing the expected squared MMD in (9) with the rational quadratic kernel [Bińkowski et al., 2018]. 404 405

406 6.1 GRAPH STRUCTURE LEARNING & OPTIMAL POINT PREDICTIONS

To test our method's ability to both calibrate the latent distribution and make optimal predictions, we train the model minimizing \mathcal{L}^{dist} as described in Section 5.2.

Figure 2 reports the validation losses during training: MMD loss \mathcal{L}^{dist} as in (6), MAE between the learned parameters θ and the ground truth θ^* as \mathcal{L}^{cal} (5), and point-prediction loss \mathcal{L}^{point} as in (4) with ℓ being the MAE. The results are averaged over 20 different model initializations and error bars report ± 1 standard deviations from the mean. Results are reported with and without applying the variance reduction (Section 5.3), by training only parameters θ while freezing ψ to ψ^* (same setting of Theorem 5.2), and by joint training of both ψ and θ .

Solving the joint learning problem Figure 2a shows that the training succeeded and the MMD 417 loss \mathcal{L}^{dist} approached its minimum (dotted line). Having minimized \mathcal{L}^{dist} , from Figure 2b we see 418 that also the calibration of latent distribution P_A^{θ} was successful; in particular, the figure shows that 419 the MAE on θ parameters $(N^{-2} \| \theta^* - \theta \|_1)$ approaches zero as training proceeds (MAE < 0.04). 420 Regarding the point predictions, Figure $\frac{2}{2c}$ confirms that \mathcal{L}^{point} reached its minimum value; recall 421 that optimal prediction MAE is not 0, as the target variable y is random, and note that a learning 422 rate reduction is applied at epoch number 5. The optimality of the point-prediction is supported also 423 by the performance on separate test data and with respect to the MSE as point-prediction loss ℓ . 424 Moreover, we observe that calibration is achieved regardless of the variance reduction and whether 425 or not parameters ψ are trained. Lastly, Figure 4 shows the learned parameters θ of the latent 426 distribution and the corresponding absolute discrepancy resulted from a (randomly chosen) training 427 run.

429 Optimization landscape of \mathcal{L}^{point} **and** \mathcal{L}^{dist} In this experiment, we analyze the values of **430** $\mathcal{L}^{point}(\psi^*, \theta)$ and $\mathcal{L}^{dist}(\psi^*, \theta)$ for different values of θ . \mathcal{L}^{point} is computed employing MAE as **431** loss function ℓ . Specifically, we let scalar p vary from 1/2 to 1 and set all $\theta_{ij} = p$ for i, j where $\theta^*_{ii} = 3/4$. Figure 1 reports the obtained results, highlighting an almost flat \mathcal{L}^{point} for values



Figure 4: The learned parameters for the latent Figure 5: Absolute error made on the parameters distribution corresponding to the stochastic adja- of the latent distribution. cency matrix.

 ≥ 0.725 . In contrast, \mathcal{L}^{dist} displays a pronounced concave shape with a clear minimum around θ^* which suggests that calibration is easier when we minimize \mathcal{L}^{dist} instead of \mathcal{L}^{point} .

450 Overall, we conclude that our approach is effective in 451 solving the joint learning problem of calibrating the latent variable while producing optimal point predictions. 452

Variance reduction effectiveness Figures 2a, 2b and 454 2c demonstrate that the proposed variance reduction 455 method (Section 5.2) yields notable advantages training 456 speed up (roughly 50% faster). For this reason, the next 457 experiments rely on variance reduction. 458

Larger graphs The theoretical results developed hold 459 for any number of nodes N. However, the number of pos-460 sible edges scales quadratically in the number of nodes. 461 In Figure 3, we show all $\sim 15K$ parameters of the con-462 sidered P_A^{θ} can be effectively learned even for relatively 463 large graphs; the final MAE on θ parameters is 0.003. 464 Note that for extremely large graphs the ratio between the 465



Figure 3: Learned θ parameters for a graph with $\sim 15K$ possible edges.

number of free parameters in θ and the size of the training set can become prohibitive. In these 466 cases, amortized learning of the edge probabilities is a potentially viable solution. 467

6.2 BEYOND ASSUMPTION 3.1

444

445

446 447

448

449

453

468

469

470 In this section, we empirically study whether Assumption 3.1 is restrictive in practical applications. 471 Specifically, we consider different degrees of model mismatch between the system model in (1) and the approximating model in (2). Unless otherwise specified, we use the same dataset and experimen-472 tal setup as described in Appendix C.1. Additional details and results are deferred to Appendix C.3. 473

474 **Perturbed** f_{ψ^*} As a first experiment, we 475 train P_A^{θ} while keeping the parameters of the 476 predictive function f_{ψ} fixed to a random pertur-477 bation of the data-generating model $f^* = f_{\psi^*}$. 478 A perturbed version of f_{ψ}^* is built by uniformly 479 drawing independent perturbation scalar values 480 $\delta_i \sim \mathcal{U}[-\Psi, \Psi]$, one for each of parameter ψ_i^* 481 of f_{ψ^*} . Then, each parameter of GNN f_{ψ} is 482 given as $\psi_i = (1 + \delta_i)\psi_i^*$. Table 1 shows that the learned latent distribution remains reason-483 ably calibrated, even when parameters can be 484 modified up to 80%. In particular, the absolute 485

Table 1: Calibration of P_A^{θ} under varying levels of misconfiguration for predictive function f_{ψ} . Results are the mean ± 1 standard deviation assessed over 8 independent runs.

Max pert. Ψ	MAE on θ	Max AE on θ
0	0.018 ± 0.005	0.12 ± 0.01
0.1	0.02 ± 0.01	0.12 ± 0.01
0.2	0.02 ± 0.02	0.14 ± 0.03
0.5	0.03 ± 0.02	0.20 ± 0.12
0.8	0.07 ± 0.02	0.36 ± 0.08

error (AE) on parameters θ is under 10% on average and increases with Ψ . Finally, Figures 8-11

show the learned parameter vectors θ for randomly extracted runs and highlight that the maximum AE of Table 1 is observed only sporadically.

489 **Generic GNN as** f_{ψ} In this second experiment, we set f_{ψ} to be a generic multilayer GNN which we jointly train with graph distribution P_A^{θ} . We comment that model family $\{f_{\psi}\}$ does not include 490 f^* , as f^* uses L-hop adjacency matrices generated from the sampled adjacency matrix A, while the 491 learnable f_{ψ} relies on multiple nonlinear 1-hop layers; details on the model architecture are reported 492 in Appendix C.3. Upon convergence, models achieved a MAE on $\theta < 0.11$ and $\mathcal{L}^{\text{point}} < 0.34$ using 493 the MAE as loss function ℓ in (4); The performance is in line with results in Figure 2c and Table 1. 494 At last, we note that because the GNN used adds self-loops, the diagonal elements of the adjacency 495 matrix are learned as zero, resulting in a larger MAE on θ (see Figure 12). However, this does not 496 impair the learning the off-diagonal θ_{ij} parameters (i.e., for $i \neq j$). Notably, in the worst-performing 497 model, these off-diagonal parameters have a MAE of 0.05. 498

Misconfigured P_A^{θ} Finally, we violate Assumption 3.1 by fixing $f_{\psi} = f^*$ and constraining some components of θ to incorrect values. Specifically, we force parameters $\theta_{i,j}$ for all edges i, j associated with nodes with id 2 and 3 in Figure 6 to 0.25, instead of the correct value of $\theta_{i,j}^* = 0.75$ as in P_A^* . Results indicate that the free parameters in θ are learned appropriately. Notably, increased uncertainty is observed for spurious edges linking to nodes in the first node community (see Figure 6). This is expected given that nearly 60% of the edges in the community were significantly downsampled. Figures 13 and 14 in Appendix C.3 show the learned parameters from randomly selected runs.

507 508

7 CONCLUSIONS

509

523 524

525

526

527

528 529

530

531 532

533

534

537

510 Graph structure learning has emerged as a research field focused on learning graph topologies in 511 support of solving downstream predictive tasks. Assuming stochastic latent graph structures, we 512 are led to a joint optimization objective: (i) learning the correct distribution of the latent topology 513 while (ii) achieving optimal predictions on the downstream task. In this paper, at first, we prove 514 both positive and negative theoretical results to demonstrate that appropriate loss functions must be 515 chosen to solve this joint learning problem. Second, we propose a sampling-based learning method that does not require the computation of the predictive likelihood. Our empirical results demonstrate 516 that this approach achieves optimal point predictions on the considered downstream task while also 517 yielding calibrated latent graph distributions. 518

Finally, we acknowledge that the proposed method requires sampling and processing multiple adjacency matrices for each input and, although the model and prediction accuracy is enhanced, a
computation overhead is requested. We plan future research to explore the applicability of this
method to other classes of neural networks beyond GNNs.

References

- Alexander Bartler, Felix Wiewel, Lukas Mauch, and Bin Yang. Training variational autoencoders with discrete latent variables using importance sampling. In 2019 27th European Signal Processing Conference (EUSIPCO), pp. 1–5. IEEE, 2019.
- James O Berger. Statistical decision theory. In *Time Series and Statistics*, pp. 277–284. Springer, 1990.

Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. In *International Conference on Learning Representations*, 2018.

- Andrea Cini, Daniele Zambon, and Cesare Alippi. Sparse graph learning from spatiotemporal time
 series. *Journal of Machine Learning Research*, 24:1–36, 2023.
- Marissa Connor, Gregory Canal, and Christopher Rozell. Variational autoencoder with learned latent structure. In *International Conference on Artificial Intelligence and Statistics*, pp. 2359–2367. PMLR, 2021.

540	Mario Coutino, Elvin Isufi, Takanori Maehara, and Geert Leus. State-space network topology iden-
541	tification from partial observations. IEEE Transactions on Signal and Information Processing
542	over Networks, 6:211-225, 2020.
543	

- Giovanni De Felice, Andrea Cini, Daniele Zambon, Vladimir Gusev, and Cesare Alippi. Graph based Virtual Sensing from Sparse and Partial Multivariate Observations. In *The Twelfth Interna- tional Conference on Learning Representations*, 2024.
- Tristan Deleu, António Góis, Chris Emezue, Mansi Rankawat, Simon Lacoste-Julien, Stefan Bauer, and Yoshua Bengio. Bayesian structure learning with generative flow networks. In *Uncertainty in Artificial Intelligence*, pp. 518–528. PMLR, 2022.
- Xiaowen Dong, Dorina Thanou, Pascal Frossard, and Pierre Vandergheynst. Learning laplacian matrix in smooth graph signal representations. *IEEE Transactions on Signal Processing*, 64(23): 6160–6173, 2016.
- Pantelis Elinas, Edwin V Bonilla, and Louis Tiao. Variational inference for graph convolutional networks in the absence of graph data and adversarial settings. Advances in Neural Information Processing Systems, 33:18648–18660, 2020.
- Bahare Fatemi, Layla El Asri, and Seyed Mehran Kazemi. Slaps: Self-supervision improves structure learning for graph neural networks. *Advances in Neural Information Processing Systems*, 34: 22667–22681, 2021.
- Matthias Fey and Jan Eric Lenssen. Fast graph representation learning with pytorch geometric. *arXiv preprint arXiv:1903.02428*, 2019.
- Alex Fout, Jonathon Byrd, Basir Shariat, and Asa Ben-Hur. Protein interface prediction using graph
 convolutional networks. *Advances in neural information processing systems*, 30, 2017.
- Luca Franceschi, Mathias Niepert, Massimiliano Pontil, and Xiao He. Learning discrete structures for graph neural networks. In *International conference on machine learning*, pp. 1972–1982.
 PMLR, 2019.
- Tilmann Gneiting. Making and Evaluating Point Forecasts. *Journal of the American Statistical Association*, 106(494):746–762, June 2011. ISSN 0162-1459. doi: 10.1198/jasa.2011.r10138.
- 571 Manuel Gomez Rodriguez, Jure Leskovec, and Bernhard Schölkopf. Structure and dynamics of
 572 information pathways in online media. In *Proceedings of the sixth ACM international conference* 573 on Web search and data mining, pp. 23–32, 2013.
- Caitlin Gray, Lewis Mitchell, and Matthew Roughan. Bayesian inference of network structure from information cascades. *IEEE Transactions on Signal and Information Processing over Networks*, 6:371–381, 2020.
- Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola.
 A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773, 2012.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.
- Gaëtan Hadjeres, Frank Nielsen, and François Pachet. Glsr-vae: Geodesic latent space regularization for variational autoencoder architectures. In 2017 IEEE symposium series on computational intelligence (SSCI), pp. 1–7. IEEE, 2017.
- Charles R Harris, K Jarrod Millman, Stéfan J Van Der Walt, Ralf Gommers, Pauli Virtanen, David
 Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J Smith, et al. Array programming with numpy. *Nature*, 585(7825):357–362, 2020.
- John D Hunter. Matplotlib: A 2d graphics environment. *Computing in science & engineering*, 9 (03):90–95, 2007.
- Bo Jiang, Ziyan Zhang, Doudou Lin, Jin Tang, and Bin Luo. Semi-supervised learning with graph
 learning-convolutional networks. In *Proceedings of the IEEE/CVF conference on computer vision* and pattern recognition, pp. 11313–11320, 2019.

594 Weonyoung Joo, Wonsung Lee, Sungrae Park, and Il-Chul Moon. Dirichlet variational autoencoder. Pattern Recognition, 107:107514, 2020. 596 Vassilis Kalofolias. How to learn a graph from smooth signals. In Artificial intelligence and statis-597 tics, pp. 920-929. PMLR, 2016. 598 Anees Kazi, Luca Cosmo, Seyed-Ahmad Ahmadi, Nassir Navab, and Michael M Bronstein. Dif-600 ferentiable graph module (dgm) for graph convolutional networks. IEEE Transactions on Pattern 601 Analysis and Machine Intelligence, 45(2):1606–1617, 2022. 602 603 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014. 604 605 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. arXiv preprint 606 arXiv:1312.6114, 2013. 607 608 Thomas Kipf, Ethan Fetaya, Kuan-Chieh Wang, Max Welling, and Richard Zemel. Neural relational 609 inference for interacting systems. In International conference on machine learning, pp. 2688– 610 2697. PMLR, 2018. 611 Volodymyr Kuleshov, Nathan Fenner, and Stefano Ermon. Accurate uncertainties for deep learning 612 using calibrated regression. In International conference on machine learning, pp. 2796–2804. 613 PMLR, 2018. 614 615 Max-Heinrich Laves, Sontje Ihler, Jacob F Fast, Lüder A Kahrs, and Tobias Ortmaier. Well-616 calibrated regression uncertainty in medical imaging with deep learning. In Medical imaging 617 with deep learning, pp. 393–412. PMLR, 2020. 618 Andrey Lokhov. Reconstructing parameters of spreading models from partial observations. Ad-619 vances in Neural Information Processing Systems, 29, 2016. 620 621 Gonzalo Mateos, Santiago Segarra, Antonio G Marques, and Alejandro Ribeiro. Connecting the 622 dots: Identifying network structure via graph signal processing. IEEE Signal Processing Maga-623 zine, 36(3):16-43, 2019. 624 Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim 625 Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement 626 learning. In International conference on machine learning, pp. 1928–1937. PMLR, 2016. 627 628 Shakir Mohamed and Balaji Lakshminarayanan. Learning in implicit generative models. 2016. 629 630 Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih. Monte carlo gradient es-631 timation in machine learning. The Journal of Machine Learning Research, 21(1):5183–5244, 2020. 632 633 Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen, Gaurav 634 Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks. 635 In Proceedings of the AAAI conference on artificial intelligence, volume 33, pp. 4602–4609, 2019. 636 637 Jishnu Mukhoti, Viveka Kulharia, Amartya Sanyal, Stuart Golodetz, Philip Torr, and Puneet Doka-638 nia. Calibrating deep neural networks using focal loss. Advances in Neural Information Process-639 ing Systems, 33:15288–15299, 2020. 640 Alfred Müller. Integral probability metrics and their generating classes of functions. Advances in 641 applied probability, 29(2):429-443, 1997. 642 643 Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? Ad-644 vances in neural information processing systems, 32, 2019. 645 Mathias Niepert, Pasquale Minervini, and Luca Franceschi. Implicit MLE: Backpropagating 646 Through Discrete Exponential Family Distributions. In Advances in Neural Information Pro-647

cessing Systems, volume 34, pp. 14567–14579. Curran Associates, Inc., 2021.

648 649 650	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. <i>Advances in neural information processing systems</i> , 32, 2019.
652 653	Xingyue Pu, Tianyue Cao, Xiaoyun Zhang, Xiaowen Dong, and Siheng Chen. Learning to learn graph topologies. <i>Advances in Neural Information Processing Systems</i> , 34:4249–4262, 2021.
654 655 656 657	Alfréd Rényi. On measures of entropy and information. In <i>Proceedings of the fourth Berkeley</i> symposium on mathematical statistics and probability, volume 1: contributions to the theory of statistics, volume 4, pp. 547–562. University of California Press, 1961.
658 659 660	Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In <i>International conference on machine learning</i> , pp. 1278–1286. PMLR, 2014.
661 662 663	Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. <i>IEEE transactions on neural networks</i> , 20(1):61–80, 2008.
664 665 666	Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. <i>Journal of Machine Learning Research</i> , 9(3), 2008.
667 668	Chao Shang, Jie Chen, and Jinbo Bi. Discrete graph structure learning for forecasting multiple time series. In <i>International Conference on Learning Representations</i> , 2021.
669 670 671	Jonathan Shlomi, Peter Battaglia, and Jean-Roch Vlimant. Graph neural networks in particle physics. <i>Machine Learning: Science and Technology</i> , 2(2):021001, 2020.
672 673 674	Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient meth- ods for reinforcement learning with function approximation. <i>Advances in neural information</i> <i>processing systems</i> , 12, 1999.
675 676 677	Gábor J Székely and Maria L Rizzo. Energy statistics: A class of statistics based on distances. <i>Journal of statistical planning and inference</i> , 143(8):1249–1272, 2013.
678 679	Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. Advances in neural information processing systems, 30, 2017.
680 681 682	Max Wasserman and Gonzalo Mateos. Graph structure learning with interpretable bayesian neural networks. <i>Transactions on machine learning research</i> , 2024.
683 684 685	Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. <i>Machine learning</i> , 8:229–256, 1992.
686 687 688	Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, and Chengqi Zhang. Graph wavenet for deep spatial-temporal graph modeling. In <i>Proceedings of the 28th International Joint Conference on Artificial Intelligence</i> , pp. 1907–1913, 2019.
689 690 691 692	Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, Xiaojun Chang, and Chengqi Zhang. Con- necting the dots: Multivariate time series forecasting with graph neural networks. In <i>Proceedings</i> <i>of the 26th ACM SIGKDD international conference on knowledge discovery & data mining</i> , pp. 753–763, 2020.
693 694 695	Jiacheng Xu and Greg Durrett. Spherical latent spaces for stable variational autoencoders. <i>arXiv</i> preprint arXiv:1808.10805, 2018.
696 697 698 699	Donghan Yu, Ruohong Zhang, Zhengbao Jiang, Yuexin Wu, and Yiming Yang. Graph-revised convolutional network. In <i>Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2020, Ghent, Belgium, September 14–18, 2020, Proceedings, Part III, pp. 378–393. Springer, 2021.</i>
100	

701 Bianca Zadrozny and Charles Elkan. Obtaining calibrated probability estimates from decision trees and naive bayesian classifiers. In *Icml*, volume 1, pp. 609–616, 2001.

Yingxue Zhang, Soumyasundar Pal, Mark Coates, and Deniz Ustebay. Bayesian graph convolutional neural networks for semi-supervised classification. In Proceedings of the AAAI conference on artificial intelligence, volume 33, pp. 5829-5836, 2019. Yanqiao Zhu, Weizhi Xu, Jinghao Zhang, Qiang Liu, Shu Wu, and Liang Wang. Deep graph struc-ture learning for robust representations: A survey. arXiv preprint arXiv:2103.03036, 14:1-1, 2021.

Appendix

A PROOFS OF THE THEORETICAL RESULTS

A.1 MINIMIZING \mathcal{L}^{point} does not guarantee calibration

In this section, we prove Proposition 4.1.

Proof of Proposition 4.1. Recall the definition of \mathcal{L}^{point} in (4) using (3)

$$\mathcal{L}^{point}(\psi,\theta) = \mathbb{E}_{x} \Big[\mathbb{E}_{y^{*} \sim P_{y|x}^{*}} \Big[\ell \big(y^{*}, T \big[P_{y|x}^{\theta,\psi} \big] \big) \Big] \Big]$$

Given loss function ℓ , T is, by definition [Berger, 1990; Gneiting, 2011], the functional that minimizes

$$\mathbb{E}_{y^* \sim P_{y|x}^*} \Big[\ell \big(y^*, T \big[P_{y|x}^* \big] \big) \Big]$$

Therefore, if $P_{y|x}^{\theta,\psi} = P_{y|x}^* \implies \mathcal{L}^{point}$ is minimal. If another distribution over y, namely, $P_{y|x}^{\psi',\theta'}$ parametrized by θ' and ψ' satisfies:

$$T\left[P_{y|x}^{\psi',\theta'}\right] = T\left[P_{y|x}^*\right]$$

almost surely on x, then,

$$\mathcal{L}^{point}(\theta',\psi') = \mathbb{E}_x \left[\mathbb{E}_{y^* \sim P_{y|x}^*} \left[\ell(y^*, T[P_{y|x}^{\psi',\theta'}]) \right] \right]$$
$$= \mathbb{E}_x \left[\mathbb{E}_{y^* \sim P_{y|x}^*} \left[\ell(y^*, T[P_{y|x}^*]) \right] \right]$$

Thus, $P_{u|x}^{\psi',\theta'}$ minimizes \mathcal{L}^{point} .

Appendix A.2 discusses graph distributions where $T[P_{y|x}^{\psi',\theta'}] = T[P_{y|x}^*]$ but $P_{y|x}^{\psi',\theta'} \neq P_{y|x}^*$. We conclude that reaching the minimum of $\mathcal{L}^{point}(\psi,\theta)$ does not imply $P_{y|x}^{\psi,\theta} = P_{y|x}^*$.

A.2 MINIMIZING \mathcal{L}^{point} does not guarantee calibration: an example with MAE

791 This section shows that \mathcal{L}^{point} equipped with MAE as ℓ admits multiple global minima for different 792 parameters θ , even for simple models and $f_{\psi} = f^*$.

Consider a single Bernoulli of parameter $\theta^* > 1/2$ as latent variable A and a scalar function $f^*(x, A)$ such that $f^*(x, 1) > f^*(x, 0)$ for all x. Given input x the value of functional $T(P_{y|x}^*)$ that minimizes

$$\mathbb{E}_{y \sim P_{y|x}^{*}} \left[\left| y - T[P_{y|x}^{*}] \right| \right] = \theta^{*} \left| f^{*}(x,1) - T[P_{y|x}^{*}] \right| + (1-\theta^{*}) \left| f^{*}(x,0) - T[P_{y|x}^{*}] \right|$$

is $T(P_{y|x}^*) = f^*(x, 1)$; this derives from the fact that range of f^* is $\{f^*(x, 0), f^*(x, 1)\}$ and the likelihood of $f^*(x, 1)$ is larger than that of $f^*(x, 0)$.

Note that $T[P_{y|x}^*] = f^*(x, 1)$ for all x, therefore also \mathcal{L}^{point} is minimized by such T. Moreover, $T[P_{y|x}^*]$ is function of θ^* and equal to $f^*(x, 1)$ for all $\theta > 1/2$. We conclude that for any $\theta \neq \theta^*$ distributions $P_{y|x}^{\theta,\psi}$ and $P_{y|x}^*$ are different, yet both of them minimize \mathcal{L}^{point} if $\theta > 1/2$.

A similar reasoning applies for $\theta^* < 1/2$.

808 A.3 MINIMIZING
$$\mathcal{L}^{dist}$$
 GUARANTEES CALIBRATION AND OPTIMAL POINT PREDICTIONS

This section proves Theorem 5.2 and a corollary of it.

Proof of Theorem 5.2. Recall from Equation (6) that

$$\mathcal{L}^{dist}(\theta) = \mathbb{E}_x \Big[\Delta(P_{y|x}^*, P_{y|x}^{\theta}) \Big]$$

815 We start by proving that if $\mathcal{L}^{dist}(\theta, \psi) = 0 \implies \mathcal{L}^{point}(\theta, \psi)$ is minimal.

Note that $\mathcal{L}^{dist}(\theta, \psi) = 0$ implies that $\Delta(P_{y|x}^*, P_{y|x}^{\theta}) = 0$ almost surely in x. Then, by Assumption 5.1, $P_{y|x}^* = P_{y|x}^{\psi,\theta}$ almost surely on x and, in particular, $T[P_{y|x}^*] = T[P_{y|x}^{\psi,\theta}]$, which leads to $\mathcal{L}^{point}(\psi, \theta)$ being minimal (Proposition 4.1).

820 We now prove that if $\mathcal{L}^{dist}(\theta, \psi^*) = 0 \implies \mathcal{L}^{cal}(\theta) = 0.$

From the previous step, we have that $\mathcal{L}^{dist}(\theta, \psi) = 0$ implies $P_{y|x}^* = P_{y|x}^{\psi,\theta}$ almost surely for $x \in I$. Under the assumption that $f_{\psi} = f_*$ and the injectivity of f_* in such $x \in I$, for any output y a single A exists such that $f_*(x, A) = y$. Therefore, the probability mass function of y equals that of A. Accordingly, $P_{y|x}^* = P_{y|x}^{\psi,\theta}$ implies $P_A^* = P_A^{\theta}$.

We also prove a corollary of Theorem 5.2.

Corollary A.1. Under Assumptions 3.1 and 5.1, if

1. $\exists \bar{x} \in Supp(P_x^*) \subseteq \mathcal{X}$ such that $f^*(\bar{x}; \cdot)$ is injective,

2. $f^*(x, A)$ is continuous in $\bar{x} \forall A \in A$,

then

$$\mathcal{L}^{dist}(\theta,\psi^*) = 0 \implies \begin{cases} \mathcal{L}^{point}(\theta,\psi^*) \text{ is minimal} \\ \mathcal{L}^{cal}(\theta) = 0. \end{cases}$$

The corollary shows that it is sufficient that f^* is continuous in x and there exists one point \bar{x} where $f^*(\bar{x}, \cdot)$ is injective to meet theorem's hypothesis $\mathbb{P}_{x \sim P_x^*}(I) > 0$; we observe that, as \mathcal{A} is discrete, the injectivity assumption is not as restrictive as if the domain were continuous.

Proof. As \mathcal{A} is a finite set, the minimum $\bar{\epsilon} = \min_{A,A' \in \mathcal{A}} \|f^*(\bar{x},A) - f^*(\bar{x},A')\| > 0$ exists and, by the injectivity assumption, is strictly positive.

By continuity of $f^*(\cdot, A)$, for every $\epsilon < \frac{1}{2}\overline{\epsilon}$ there exists δ , such that for all $x \in B(\overline{x}, \delta)$ we have $||f^*(\overline{x}, A) - f^*(x, A)|| < \epsilon$. It follows that, $\forall x \in B$,

 $\begin{aligned} \|f^*(x,A) - f^*(x,A')\| \\ &\geq \|f^*(\bar{x},A) - f^*(\bar{x},A')\| - \|f^*(\bar{x},A) - f^*(x,A)\| - \|f^*(\bar{x},A') - f^*(x,A')\| \\ &\geq \|f^*(\bar{x},A) - f^*(\bar{x},A')\| - 2\epsilon \\ &\geq \|f^*(\bar{x},A) - f^*(\bar{x},A')\| - \bar{\epsilon} > 0. \end{aligned}$

Finally, as $\bar{x} \in \text{Supp}(P_x^*)$ and $B(\bar{x}, \delta) \subseteq I$, we conclude that

$$\mathbb{P}_x(I) \ge \mathbb{P}_x(B(\bar{x},\delta)) > 0.$$

therefore, we are in the hypothesis of Theorem 5.2 and can conclude that

$$\mathcal{L}^{dist}(\theta, \psi^*) = 0 \implies \begin{cases} \mathcal{L}^{point}(\theta, \psi^*) \text{ is minimal} \\ \mathcal{L}^{cal}(\theta) = 0. \end{cases}$$

864 A.4 INJECTIVITY HYPOTHESIS FOR GRAPH NEURAL NETWORKS 865

866 Now, we show that hypothesis $\mathbb{P}_{x \sim P_x^*}(I) > 0$ of Theorem 5.2 is always met for certain families of graph neural networks. 867

868 **Proposition A.2.** Consider a 1-layer GNN of the form $f^*(x, A) : \sigma(Ax) = y$, with $x, y \in \mathbb{R}^N$ and nonlinear bijective activation function σ . If the support $Supp(P_x^*)$ of x contains any ball B in \mathbb{R}^N 870 then $\mathbb{P}_{x \sim P_{\pi}^*}(I) > 0.$

To prove Proposition A.2, we rely on following lemma. 872

Lemma A.3. Given g(x,a) = ax with $a \in \{0,1\}^{1 \times N}$ and $x \in \mathbb{R}^{N \times 1}$. Let $I_g = \{x : x \in \mathbb{R}^{N \times 1}\}$ 873 g(x,a) is injective in $a \subseteq \mathcal{X}$ be the set of points $x \in \mathcal{X}$ such that map $a \mapsto g(x,a)$ is injective. 874 The following implication holds: 875

 $x \notin I_q \iff \exists \delta \neq \mathbf{0} \in \{-1, 0, 1\}^{1 \times N}$ s.t. $\delta \perp x$. (14)

Proof. We prove the two implications separately.

- (\Longrightarrow) If $x \notin I_g$, then there exist $a', a'' \in \{0, 1\}^{1 \times N}$ with $a' \neq a''$ such that a'x = a''x. This implies that (a' a'')x = 0. Defining δ as (a' a''), we have proven that there exist $\delta \neq \mathbf{0} \in \{-1, 0, 1\}^{1 \times N}$ such that $\delta x = 0$, i.e., $\delta \perp x$.
- (\Leftarrow) Assume that $\exists \delta \neq \mathbf{0} \in \{-1, 0, 1\}^{1 \times N}$ such that $\delta \perp x$. Each component δ_i of δ can be written as the difference between two values $a'_i, a''_i \in \{0, 1\}$. As $\delta \neq 0$ then there exists at least one index $j \in \{1, ..., N\}$ such that $a'_j \neq a''_j$. This implies that $\exists a', a'' \in \{0, 1\}^{1 \times N}$ with $a' \neq a''$ s.t. (a' - a'')x = 0, which implies that $x \notin I_q$.

Proof of Proposition A.2. We begin by considering the projection $\bar{g}(x, a) = ax$ with $a \in \{0, 1\}^{1 \times N}$ and $x \in \mathbb{R}^N$. Then we extend to $A \in \{0, 1\}^{N \times N}$ and to nonlinear functions. 890 891

892 Let $I_{\bar{g}}^C = \mathbb{R}^N \setminus I_{\bar{g}}$ be the complement in \mathbb{R}^N of $I_{\bar{g}}$. Recalling Lemma A.3 and its notation, we have $3^N - 1$ possible δ , defining a collection of $(3^N - 1)/2$ hyperplanes of vectors x perpendicular 893 894 to at least one δ ; set $I_{\bar{g}}^C$ is the union of such a finite collection of hyperplanes. By hypothesis, 895 $\operatorname{Supp}(P_x^*)$ contains a ball $B \in \mathbb{R}^N$, therfore $\operatorname{Supp}(P_x^*) \not\subset I_{\overline{q}}^C$ and $\mathbb{P}_{x \sim P_x^*}(I_{\overline{q}}^C) < 1$. We conclude 896 that $\mathbb{P}_{x \sim P_x^*}(I_{\bar{q}}) = 1 - \mathbb{P}_{x \sim P_x^*}(I_{\bar{q}}^C) > 0.$ 897

898 A similar result is proven for $\overline{G}(x,A) = Ax$ with $A \in \{0,1\}^{N \times N}$. In fact, \overline{G} is a stack of N 899 functions \bar{g} above and $I_{\bar{G}} = I_{\bar{g}}$. Finally, composing injective function G with injective function σ leads to function $g(x, A) = \sigma(G(x, A))$ being injective in A for the same points x for which G is 900 901 injective, thus proving the proposition.

В ESTIMATION OF OPTIMAL β_1 AND β_2

Here we show that, when reducing the variance of the SFE via control variates in (12), the best β_1 and β_2 can be approximated by

$$\tilde{\beta}_{1} = \mathbb{E}_{\substack{x \sim P_{x}^{*} \\ A_{1}, A_{2} \sim P_{A}^{\theta}}} \left[\kappa \left(f_{\psi}(x, A_{1}), f_{\psi}(x, A_{2}) \right) \right], \qquad \tilde{\beta}_{2} = \mathbb{E}_{\substack{(x, y^{*}) \sim P_{x, y}^{*} \\ A \sim P_{A}^{\theta}}} \left[\kappa \left(y^{*}, f_{\psi}(x, A) \right) \right], \qquad (15)$$

911 Consider generic function L(A) depending on a sample A of a parametric distribution $P_A^{\theta}(A)$ and the surrogate loss $\tilde{L}(A)$ in (11), i.e., 913

$$\tilde{L}(A) = L(A) - \beta \Big(h(A) - \mathbb{E}_{A \sim P^{\theta}}[h(A)] \Big);$$
(16)

915 This choice is not new in the literature [Sutton et al., 1999; Mnih et al., 2016] where β is often 916 referred to as baseline. The 1-sample MC approximation of the loss becomes 917

$$\nabla_{\theta} \mathbb{E}_{A \sim P^{\theta}} [L(A)] \approx \tilde{L}(A') \nabla_{\theta} \log P^{\theta}(A') = (L(A') - \beta) \nabla_{\theta} \log P^{\theta}(A'), \tag{17}$$

912

914

902 903

904 905

906

871

876 877 878

879

884

885 886

887 888



Figure 6: The adjacency matrices used in this paper are sampled from this graph. Each edge in orange is independently sampled with probability θ^* . In the picture, 3 communities of an arbitrarily large graph are shown.

with A' sampled from P_A^{θ} . The variance of the estimator is

$$\mathbb{V}_{A\sim P^{\theta}}\left[(L(A)-\beta)\nabla_{\theta}\log P^{\theta}(A)\right] = \mathbb{V}_{A\sim P^{\theta}}\left[L(A)\nabla_{\theta}\log P^{\theta}(A)\right] + \beta^{2} \mathbb{E}_{A\sim P^{\theta}}\left[\left(\nabla_{\theta}\log P^{\theta}(A)\right)^{2}\right] - 2\beta \mathbb{E}_{A\sim P^{\theta}}\left[L(A)\left(\nabla_{\theta}\log P^{\theta}(A)\right)^{2}\right]$$
(18)

and the optimal value β that minimizes it is

$$\tilde{\beta} = \frac{\mathbb{E}_{A \sim P^{\theta}} \left[L(A) \left(\nabla_{\theta} \log P^{\theta}(A) \right)^{2} \right]}{\mathbb{E}_{A \sim P^{\theta}} \left[\left(\nabla_{\theta} \log P^{\theta}(A) \right)^{2} \right]}.$$
(19)

If we approximate the numerator with $\mathbb{E}[L(A)]\mathbb{E}[(\nabla_{\theta} \log P^{\theta}(A))^2]$, we obtain that $\tilde{\beta} \approx \mathbb{E}[L(A)]$. By substituting L(A) with the two terms of (10) we get the values of β_1 and β_2 in (15).

We experimentally validate the effectiveness of this choice of β in Section 6.

C FURTHER EXPERIMENTAL DETAILS

C.1 DATASET DESCRIPTION AND MODELS

In this section, we describe the considered synthetic dataset, generated from the system model (1). The latent graph distribution P_A^* is a multivariate Bernoulli distribution of parameters θ_{ij}^* : $P_A^* \equiv P_{\theta^*}(A) = \prod_{ij} \theta_{ij}^{*A_{ij}} (1 - \theta_{ij}^*)^{(1 - A_{ij})}$. The components of θ^* are all null, except for the edges of the graph depicted in Figure 6 which are set to 3/4. A heatmap of the adjacency matrix can be found in Figure 7.







Table 2: Table of the parameters used to generate the synthetic dataset.

θ^*	0.75
σ_x	1.5
N	12
d_{in}	4
d_{out}	1
ψ_1^*	[-0.2, 0.4, -0.8, 0.6]
ψ_2^*	[-0.3, 0.8, 0.2, -0.7]

Regarding the GNN function f^* , we use the following system model:

$$\begin{cases} y = f_{\psi^*}(A, x) = \tanh\left(\sum_{l=1}^L \mathbb{1}[A^l \neq 0] x \psi_l^*\right) \\ A \sim P_{\theta^*}(A) \end{cases}$$
(20)

where $\mathbb{1}[\cdot]$ is the element-wise indicator function: $\mathbb{1}[a] = 1 \iff a$ is true. $x \in \mathbb{R}^{N \times d_{in}}$ are randomly generated inputs: $x \sim \mathcal{N}(0, \sigma_x^2 \mathbb{I})$. $\psi_l^* \in \mathbb{R}^{d_{out} \times d_{in}}$ are part of the system model parameters. We summarize the parameters considered in our experiment in Table 2.

The approximating model family (2) used in the experiment is the same as the data-generating process, with all components of parameter vectors θ and ψ being trainable. The squared MMD discrepancy is defined over Rational Quadratic kernel [Bińkowski et al., 2018]

$$\kappa(y',y'') = \left(1 + \frac{\|y' - y''\|_2^2}{2\,\alpha\,\sigma^2}\right)^{-\alpha} \tag{21}$$

1000 of parameters $\sigma = 0.7$ and $\alpha = 0.02$.

1002 The model is trained using Adam optimizer [Kingma & Ba, 2014] with parameters $\beta_1 = 0.6$, 1003 $\beta_2 = 0.95$. Where not specified, the learning rate is set to 0.1 and decreased to 0.01 after 5 epochs. 1004 We grouped data points into batches of size 128. Initial values of θ are independently sampled from 1005 the $\mathcal{U}(0.25, 0.35)$ uniform distribution.

1006

1008

972

992 993

994

995

996

997 998 999

7 C.2 DESCRIPTION OF THE EXPERIMENT IN SECTION 4

In this experiment, we generate 512 data points using the system model described in Appendix C.1. We construct a model identical to the system model, except that $\theta_{ij} = p$ for all i, j where $\theta_{i,j}^* = 0.75$ and 0 elsewhere. We vary scalar p from 0.5 to 1 with steps of 0.025. Therefore, only the model with p = 0.75 is identical to the data-generating model.

For each input x in the dataset, a point prediction is produced by sampling $N_{adj} = 32$ adjacency matrices and computing the median. This approach allows to estimate \mathcal{L}^{point} using the MAE as loss function ℓ , as depicted by the red points in Figure 1, for different values of θ . For comparison purposes, we estimate \mathcal{L}^{dist} using the maximum mean discrepancy as proposed in Section 5.

- 1017
- 1018 C.3 ADDITIONAL DETAILS OF SECTION 6.2

We present here additional Figures discussed in Section 6.2.

Fixed perturbed f_{ψ} Figures in this paragraph correspond to the experiment where the processing function f_{ψ} is fixed on a perturbed version of f^* . Figures 8 – 11 correspond to runs with increasing perturbation factor Ψ .



Figure 8: Learned θ_{ij} parameters (a) and Absolute Error (b) for maximum perturbation factor Ψ of 10%.

Figure 9: Learned θ_{ij} parameters (a) and Absolute Error (b) for maximum perturbation factor Ψ of 20%.



Figure 10: Learned θ_{ij} parameters (a) and Absolute Error (b) for maximum perturbation factor Ψ of 50%.

Figure 11: Learned θ_{ij} parameters (a) and Absolute Error (b) for maximum perturbation factor Ψ of 80%.

1055 Generic GNN as f_{ψ} To evaluate our ap-1056 proach in a more realistic setting, we use a 1057 generic GNN as f_{ψ} . Specifically, we imple-1058 ment GNNs from [Morris et al., 2019] with 1059 varying numbers of layers and layer sizes. It is important to note that the GNN implementa-1061 tion includes self-loops, which prevents the diagonal elements from being correctly learned. 1062 However, this does not impede our method 1063 from learning the remaining edges accurately. 1064

Table 3 presents the network configurations and whether they successfully converged to the ground truth distribution. Since diagonal elements artificially inflate the MAE for θ , we consider a model to have converged if the final MAE on θ is less than 0.11.

Table 3: Network configurations and corresponding convergence results.

Layers dimensions	Convergence
[4, 1]	x
[4, 1, 1]	x
[4, 2, 1]	\checkmark
[4, 8, 1]	\checkmark
[4, 8, 2, 1]	\checkmark
[4, 16, 8, 1]	\checkmark
[4, 32, 8, 1]	\checkmark
[4, 64, 8, 1]	\checkmark
[4, 64, 16, 1]	\checkmark
[4, 64, 32, 1]	\checkmark
$\left[4, 8, 8, 4, 1 ight]$	\checkmark

1071 Most of the models successfully converged, except those with high bias. This demonstrates that our 1072 method is effective even beyond Assumption 3.1. In Figure 12 we show the learned parameters of 1073 P_A^{θ} for a randomly extracted run.

1074

1036

1037

1038

1050

1051

1052 1053 1054



1076 **Misconfigured** P_A^{θ} Figures 13 and 14 correspond to the experiment where some θ_{ij} values of P_A^{θ} 1077 are fixed at incorrect values, while the processing function f_{ψ} is fixed to the true one. In the commu-1078 nity affected by the perturbation, free θ_{ij} values tend to be sampled more frequently to compensate 1079 for the downsampling imposed by the perturbation. Interestingly, all the edges with at least one edge 1079 in the second community (75% of the edges) appear unaffected by the perturbation.



Figure 12: (a) Learned θ_{ij} parameters when the parametric processing function f_{ψ} is a generic GNN as presented in [Morris et al., 2019] and (b) Absolute Error made with respect to true parameters θ_{ii}^* . As self-loops are deterministically added by the network, the diagonal elements should not be considered.



Figure 13: Learned θ_{ij} parameters (a) and Abso-Figure 14: Learned θ_{ij} parameters (a) and Abso-lute Error (b) for misconfigured P_A^{θ}

lute Error (b) for misconfigured P_A^{θ}

C.4 COMPUTE RESOURCES AND OPEN-SOURCE SOFTWARE

The paper's experiments were run on a workstation with AMD EPYC 7513 processors and NVIDIA RTX A5000 GPUs; on average, a single model training terminates in a few minutes with a memory usage of about 2GB.

The developed code relies on PyTorch [Paszke et al., 2019] and the following additional open-source libraries: PyTorch Geometric [Fey & Lenssen, 2019], NumPy [Harris et al., 2020] and Mat-plotlib [Hunter, 2007].