

Appendix

In the appendix, we first show the ELBOs for the CVAEs under acyclic correlation graphs, as well as the derivation for the ELBO of CVAE_{corr} (Eq. 8) for general correlation graphs. Also we show the counter example mentioned in Section 3.1. Then, we show the algorithm details for optimization with the loss function in Eq. 9 and the proofs of Theorems 2 and 3. At the end, we provide more information related to the models and inference settings.

A. ELBOs and Derivations

ELBO for CVAE_{ind} with acyclic graphs. CVAE_{ind} applies a correlated prior $p_0^{\text{corr}}(\mathbf{z})$ as in Eq. 4. Therefore,

$$\begin{aligned} & \mathcal{L}^{\text{CVAE}_{\text{ind-acyclic}}(\boldsymbol{\lambda}, \boldsymbol{\theta})} \\ &= \mathbb{E}_{q_{\boldsymbol{\lambda}}(\mathbf{z}|\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}|\mathbf{x})||p_0(\mathbf{z})) \\ &= \sum_{i=1}^n \left(\mathbb{E}_{q_{\boldsymbol{\lambda}}(\mathbf{z}_i|\mathbf{x}_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i|\mathbf{z}_i)] - \text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_i|\mathbf{x}_i)||p_0(\mathbf{z}_i)) \right) \\ & \quad + \sum_{(v_i, v_j) \in E} \mathbb{E}_{q_{\boldsymbol{\lambda}}(\mathbf{z}_i|\mathbf{x}_i)q_{\boldsymbol{\lambda}}(\mathbf{z}_j|\mathbf{x}_j)} \log \frac{p_0(\mathbf{z}_i, \mathbf{z}_j)}{p_0(\mathbf{z}_i)p_0(\mathbf{z}_j)}. \end{aligned} \quad (10)$$

This ELBO is a lower bound of the log probability $p(\mathbf{x}; \boldsymbol{\theta})$ under the correlated prior (as in Eq. 4). Optimizing this ELBO helps the variational distributions to regularize according to the correlation graph G .

ELBO for CVAE_{corr} with acyclic graphs. By changing q to be the correlated distribution as in Eq. 5, the ELBO in Eq. 10 becomes

$$\begin{aligned} & \mathcal{L}^{\text{CVAE}_{\text{corr-acyclic}}(\boldsymbol{\lambda}, \boldsymbol{\theta})} \\ &= \sum_{i=1}^n \left(\mathbb{E}_{q_{\boldsymbol{\lambda}}(\mathbf{z}_i|\mathbf{x}_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i|\mathbf{z}_i)] - \text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_i|\mathbf{x}_i)||p_0(\mathbf{z}_i)) \right) \\ & \quad - \sum_{(v_i, v_j) \in E} \left(\text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_i, \mathbf{z}_j|\mathbf{x}_i, \mathbf{x}_j)||p_0(\mathbf{z}_i, \mathbf{z}_j)) \right. \\ & \quad \left. - \text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_i|\mathbf{x}_i)||p_0(\mathbf{z}_i)) - \text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_j|\mathbf{x}_j)||p_0(\mathbf{z}_j)) \right). \end{aligned} \quad (11)$$

Derivation for Eq. 8. We derive the ELBO of CVAE_{corr} (Eq. 8) for general correlation graphs here. Recall that Eq. 7 shows that

$$\begin{aligned} \log p_{\boldsymbol{\theta}}(\mathbf{x}) &\geq \frac{1}{|\mathcal{A}_G|} \sum_{G' \in \mathcal{A}_G} \left(\mathbb{E}_{q_{\boldsymbol{\lambda}}^{G'}(\mathbf{z}|\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] \right. \\ & \quad \left. - \text{KL}(q_{\boldsymbol{\lambda}}^{G'}(\mathbf{z}|\mathbf{x})||p_0^{G'}(\mathbf{z})) \right). \end{aligned}$$

By the definition of \mathcal{A}_G , the right hand side of the above inequality equals the following sum

$$\begin{aligned} & \sum_{i=1}^n \left(\mathbb{E}_{q_{\boldsymbol{\lambda}}(\mathbf{z}_i|\mathbf{x}_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i|\mathbf{z}_i)] - \text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_i|\mathbf{x}_i)||p_0(\mathbf{z}_i)) \right) \\ & \quad - \frac{1}{|\mathcal{A}_G|} \sum_{G' \in \mathcal{A}_G} \sum_{(v_i, v_j) \in E'} \left(\text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_i, \mathbf{z}_j|\mathbf{x}_i, \mathbf{x}_j)||p_0(\mathbf{z}_i, \mathbf{z}_j)) \right. \\ & \quad \left. - \text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_i|\mathbf{x}_i)||p_0(\mathbf{z}_i)) - \text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_j|\mathbf{x}_j)||p_0(\mathbf{z}_j)) \right) \end{aligned}$$

The pairwise sum part of the above equation is an average over a sum over all edges of all maximal acyclic subgraphs of G . Therefore, for each edge $e = (v_i, v_j) \in E$, the number of times it appears in this pairwise sum part of the above sum is just the number of maximal acyclic subgraphs containing this edge. Therefore, this part can be viewed as a weighed sum over all edges in E , where the weights come from the fraction ratios in Definition 2. With this definition, we can further write the above sum as

$$\begin{aligned} & \sum_{i=1}^n \left(\mathbb{E}_{q_{\boldsymbol{\lambda}}(\mathbf{z}_i|\mathbf{x}_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i|\mathbf{z}_i)] - \text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_i|\mathbf{x}_i)||p_0(\mathbf{z}_i)) \right) \\ & \quad - \sum_{(v_i, v_j) \in E} w_{G, (v_i, v_j)}^{\text{MAS}} \left(\text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_i, \mathbf{z}_j|\mathbf{x}_i, \mathbf{x}_j)||p_0(\mathbf{z}_i, \mathbf{z}_j)) \right. \\ & \quad \left. - \text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_i|\mathbf{x}_i)||p_0(\mathbf{z}_i)) - \text{KL}(q_{\boldsymbol{\lambda}}(\mathbf{z}_j|\mathbf{x}_j)||p_0(\mathbf{z}_j)) \right) \\ & = \mathcal{L}^{\text{CVAE}_{\text{corr}}(\boldsymbol{\lambda}, \boldsymbol{\theta})}, \end{aligned}$$

which is exactly Eq. 8.

B. Counter example for Section 3.1

In Section 3.1, we mentioned that directly applying the ELBOs we derived (equations shown in Appendix A) will fail partly due to that optimizing over these ELBOs may lead the algorithm to learn useless parameters that lead the loss goes to infinity due to that these two equations are not guaranteed to be a lower bound of some log-likelihood function for general graphs. Here we provide a simple example to illustrate this.

Let us consider $G = (V, E)$ is a K_4 complete graph with $|V| = 4$ vertices and $\binom{|V|}{2} = 6$ edges. For simplicity, we consider the latent variables $z_1, z_2, z_3, z_4 \in \mathbb{R}$ and the prior distribution $p_0(z_i, z_j) = \mathcal{N}\left(\begin{pmatrix} z_i \\ z_j \end{pmatrix}; \boldsymbol{\mu} = \mathbf{0}_2, \Sigma = \begin{pmatrix} 1 & \tau \\ \tau & 1 \end{pmatrix}\right)$ for some $\tau \in (0, 1)$. If we extend the CVAE_{ind} and for the variational distribution, we set $q(z_i|\mathbf{x}_i)$ to be the normal distribution $N(\mu_i, \sigma_i^2)$, by

simply apply Eq. 10, then the loss function becomes

$$\mathcal{L} = \sum_{i=1}^4 \left(\mathbb{E}_{q_{\lambda}(z_i|\mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i|z_i)] + \mu_i^2 - \frac{1 + 2\tau^2}{2(1 - \tau^2)} \sigma_i^2 + \ln(\sigma_i) \right) - \frac{1}{2(1 - \tau^2)} \sum_{1 \leq i < j \leq 4} (\mu_i^2 + \mu_j^2 - 2\tau\mu_i\mu_j).$$

If we maintain σ_i 's unchanged, set the model parameter θ in a way that makes $p_{\theta}(\mathbf{x}_i|z_i)$ unrelated to z_i (e.g. set the parameter that multiply with z_i to be 0) and set $\mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu$ and let $\mu \rightarrow \infty$, then \mathcal{L} will go to $+\infty$ if $\tau > \frac{1}{2}$. Therefore, directly applying Eq. 10 does not work. Directly applying Eq. 11 (i.e. extending the CVAE_{corr}) will make the result even worse since it always has an optimal value at least as high as Eq. 10. In general, any general graphs with K_4 subgraph may suffer from the issue we just mentioned. We can not obtain useful latent embeddings by directly applying the Eqs. 10 and 11.

C. Algorithm details for optimization with the loss $\mathcal{L}^{\text{CVAE}_{\text{corr-NS}}(\lambda, \theta)}$

We show the details of optimization with the loss $\mathcal{L}^{\text{CVAE}_{\text{corr-NS}}(\lambda, \theta)}$ (Eq. 9) as in Algorithm 2.

Algorithm 2 Optimization with the loss $\mathcal{L}^{\text{CVAE}_{\text{corr-NS}}(\lambda, \theta)}$

Input: undirected graph $G = (V = \{v_1, \dots, v_n\}, E)$, input data $\mathbf{x}_1, \dots, \mathbf{x}_n$, the parameter $\gamma > 0$.

Compute the edge weights $w_{G,e}^{\text{MAS}}$ using Algorithm 1.

Initialize the parameters λ, θ .

while Not converged **do**

 Compute the gradient $\nabla_{\lambda, \theta} \mathcal{L}^{\text{CVAE}_{\text{corr-NS}}(\lambda, \theta)}$.

 Update the parameters λ and θ using this gradient.

end while

Return The parameters λ, θ .

If we subsample the vertices in V for the singleton part of this loss, subsample edges in E for the ‘‘positive’’ sample pairwise part of this loss and subsample edges of the complete graph K_n for the ‘‘negative’’ sample pairwise part of this loss, then we get the stochastic optimization version for this algorithm.

D. Proofs

We prove the Theorems 2 and 3 here.

D.1. Proof of Theorem 2

Proof. Given the graph $G = (V, E)$ and the edge $(v_i, v_j) \in E$. Denote G 's Laplacian matrix as L . Also denote $L_{-a,-b}$ as the sub-matrix of L after deleting the a^{th} row and the

b^{th} column. Denote $L_{-ab,-cd}$ as the sub-matrix of L after deleting the $a^{\text{th}}, b^{\text{th}}$ rows and the $c^{\text{th}}, d^{\text{th}}$ columns.

By Matrix Tree Theorem (Theorem 1 (Chaiken & Kleitman, 1978)), we know that the (i, i) -cofactor $C_{i,i} = |L_{-i,-i}|$ is the number of spanning trees of G .

Construct a graph $G' = (V, E/\{v_i, v_j\})$, i.e. the graph G after removing the edge (v_i, v_j) . Denote the Laplacian matrix of G' as L' . Then we will find that the matrix $L'_{-i,-i}$ is the same with $L_{-i,-i}$ except that they $L_{-i,-i}$ is 1 larger than $L'_{-i,-i}$ on the entry at (j, j) . By Matrix Tree Theorem, $|L'_{-i,-i}|$ is the number of spanning trees of G' . Since G differs from G' by only having one more edge (v_i, v_j) , we know that $|L_{-i,-i}| - |L'_{-i,-i}|$ represents the number of spanning trees in G that contains the edge (v_i, v_j) .

Denote the entry at (j, k) of $L_{-i,-i}$ as $L_{-i,-i;j,k}$. Since we know that

$$\begin{cases} |L_{-i,-i}| &= \sum_{k \neq i} (-1)^{j+k} L_{-i,-i;j,k} |L_{-ij,-ik}|, \\ |L'_{-i,-i}| &= \sum_{k \neq i} (-1)^{j+k} L'_{-i,-i;j,k} |L'_{-ij,-ik}|. \end{cases}$$

Subtract the second equation from the first one we get

$$|L_{-i,-i}| - |L'_{-i,-i}| = (-1)^{2j} |L_{-ij,-ij}|$$

which is just the complement of the Minor $M_{ij,ij}$ of L . Hence, the number of spanning trees of G that contains the edge (v_i, v_j) is $M_{ij,ij}$, the determinant of the sub-matrix of the Laplacian matrix L of G , after deleting the the $i^{\text{th}}, j^{\text{th}}$ rows and the $i^{\text{th}}, j^{\text{th}}$ columns. \square

D.2. Proof of Theorem 3

Proof. We borrow the notations from Appendix D.1. Given the undirected connected graph $G = (V, E)$ and an edge $(v_i, v_j) \in E$, we want to compute the ratio $w_{G,(v_i,v_j)}^{\text{MAS}}$. Since G is connected, this ratio is just the fraction of G 's spanning trees containing the edge (v_i, v_j) . By Theorem 1 and 2, this ratio is just $\frac{|L_{-ij,-ij}|}{|L_{-i,-i}|}$.

Since G is connected, it contains at least one spanning tree. Hence $|L_{-i,-i}| > 0$, which means $L_{-i,-i}$ is invertible. Therefore, we know that $\frac{|L_{-ij,-ij}|}{|L_{-i,-i}|} = L_{-i,-i;j,j}^{-1}$.

Consider the original Laplacian matrix L before deleting any row and column. Denote $|V| = n$. Since L is always symmetric and always have an eigenvector $\mathbf{v}_n = \frac{1}{\sqrt{n}} \mathbf{1}_n$ with corresponding eigenvalue $\lambda_n = 0$, we perform eigenvalue decomposition on L and write L as:

$$L = \sum_{k=1}^n \lambda_k \mathbf{v}_k \mathbf{v}_k^{\top} = \sum_{k=1}^{n-1} \lambda_k \mathbf{v}_k \mathbf{v}_k^{\top}.$$

Where $\lambda_1, \dots, \lambda_{n-1}, \lambda_n$ are L 's eigenvalues and $\mathbf{v}_1, \dots, \mathbf{v}_{n-1}, \mathbf{v}_n$ are the corresponding orthogonal

unit eigenvectors (i.e. $Q_v = (\mathbf{v}_1 \ \cdots \ \mathbf{v}_{n-1} \ \mathbf{v}_n)$ is an orthogonal matrix).

Denote $v_{a,b}$ as the b -th coordinate of \mathbf{v}_a and $\mathbf{v}_{a,-b}$ as the sub-vector \mathbf{v}_a after deleting the b^{th} coordinate. Then

$$L_{-i,-i} = \sum_{k=1}^{n-1} \lambda_k \mathbf{v}_{k,-i} \mathbf{v}_{k,-i}^\top.$$

$L_{-i,-i}$ is invertible, which is of rank $n-1$. While each of the matrix $\mathbf{v}_{k,-i} \mathbf{v}_{k,-i}^\top$ is of rank 1. Hence, we must have $\lambda_i \neq 0$ for all $i \in \{1, \dots, n-1\}$.

Construct vectors $\mathbf{u}_1, \dots, \mathbf{u}_n \in \mathbb{R}^n$ such that

$$\mathbf{u}_{k,a} = \begin{cases} v_{k,a} - v_{k,i} & \text{if } a \neq i \\ v_{k,a} & \text{if } a = i. \end{cases} \quad (12)$$

Also, construct the matrix

$$U = \sum_{k=1}^{n-1} \lambda_k^{-1} \mathbf{u}_k \mathbf{u}_k^\top.$$

Since we know that $(\mathbf{v}_1 \ \cdots \ \mathbf{v}_{n-1} \ \mathbf{v}_n)$ forms an orthogonal basis and $\mathbf{v}_n = \frac{1}{\sqrt{n}} \mathbf{1}_n$, it is easy to see (after simple calculations) that

$$\mathbf{u}_{k,-i}^\top \cdot \mathbf{v}_{k',-i} = \begin{cases} 1 & \text{if } k = k' \\ 0 & \text{if } k \neq k'. \end{cases}$$

Therefore, we will have

$$U_{-i,-i} L_{-i,-i} = I_{n-1}$$

which indicates that $U_{-i,-i} = L_{-i,-i}^{-1}$. Hence, the ratio we want to find is just $U_{-i,-i;j,j}$.

Recall the definition of $\mathbf{u}_1, \dots, \mathbf{u}_n$ in Eq. 12, denote $Q_u = (\mathbf{u}_1 \ \cdots \ \mathbf{u}_{n-1} \ \mathbf{u}_n)$, we get $Q_u = P_i Q_v$, where

$$P_i = \begin{pmatrix} 1 & & & -1 & & & & & \\ & 1 & & -1 & & & & & \\ & & \ddots & \vdots & & & & & \\ & & & 1 & & & & & \\ & & & -1 & 1 & & & & \\ & & & \vdots & & \ddots & & & \\ & & & -1 & & & 1 & & \\ & & & -1 & & & & & 1 \end{pmatrix} \in \mathbb{R}^{n \times n}$$

where the -1 's appear on the i^{th} column. Hence, denote $D = \text{diag}(\lambda_1, \dots, \lambda_n)$, we have

$$U = Q_u D^+ Q_u^\top = P_i Q_v D^+ Q_v^\top P_i^\top = P_i L^+ P_i^\top.$$

Therefore, the ratio $w_{G,(v_i,v_j)}^{\text{MAS}}$ we want to find, which is equal to $U_{-i,-i;j,j}$, is just the (j,j) -entry of $P_i L^+ P_i^\top$, which is

$$L_{i,i}^+ - L_{i,j}^+ - L_{j,i}^+ + L_{j,j}^+.$$

□

E. Additional experiment settings

We provide some additional information related to the models and inference settings here.

For all methods, we set the latent embeddings to have the dimensionality $d = 100$.

For the standard VAEs, the CVAE_{ind} and the CVAE_{corr}, we apply a two-layer feed-forward neural network for the generative model $p_\theta(\mathbf{x}_i | \mathbf{z}_i)$ and a two-layer feed-forward neural network for the singleton variational approximations $q_\lambda(\mathbf{z}_i | \mathbf{x}_i)$. The model likelihood functions $p_\theta(\mathbf{x} | \mathbf{z})$'s are Multinomial distribution (for the user matching and link prediction experiments) and Bernoulli distribution (for the spectral clustering experiments). The singleton variational approximations $q_\lambda(\mathbf{z}_i | \mathbf{x}_i)$ are all diagonal normal distributions. The singleton prior density function $p_0(\cdot)$ is the standard normal distribution.

For the CVAEs, we set the pairwise prior density function $p_0(\cdot) = \mathcal{N}\left(\boldsymbol{\mu} = \mathbf{0}_{2d}, \Sigma = \begin{pmatrix} I_d & \tau \cdot I_d \\ \tau \cdot I_d & I_d \end{pmatrix}\right)$ for $\tau = 0.99$. It can be seen that, the singleton prior density function $p_0(\cdot)$ and the pairwise prior density function $p_0(\cdot, \cdot)$ satisfy the constraints in Eq. 2. For the CVAE_{corr}, we treat $q_\lambda(\mathbf{z}_i, \mathbf{z}_j | \mathbf{x}_i, \mathbf{x}_j)$ as a multivariate normal distributions such that the covariance matrices $\text{Cov}(\mathbf{z}_i, \mathbf{z}_i)$, $\text{Cov}(\mathbf{z}_j, \mathbf{z}_j)$ and $\text{Cov}(\mathbf{z}_i, \mathbf{z}_j)$ are all diagonal matrices. Instead of only learning the singleton density function $q_\lambda(\cdot | \cdot)$, the CVAE_{corr} also learn a two-layer feedforward neural network that takes the concatenation $(\mathbf{x}_i, \mathbf{x}_j)$ as input and output the covariance between \mathbf{z}_i and \mathbf{z}_j on each of these d dimensions. As a result, $q_\lambda(\mathbf{z}_i, \mathbf{z}_j | \mathbf{x}_i, \mathbf{x}_j)$ can be factorized as a product of d bi-variate normal distributions, whose marginal distributions on \mathbf{z}_i and \mathbf{z}_j are consistent with the singleton variational approximations $q_\lambda(\mathbf{z}_i | \mathbf{x}_i)$ and $q_\lambda(\mathbf{z}_j | \mathbf{x}_j)$, respectively.

For GraphSAGE, we choose to use $K = 2$ aggregation steps and use the mean aggregator function. We use $Q = 20$ negative samples to optimize the loss function.

For all methods, we apply stochastic gradient optimizations with a step size of 10^{-3} . We use the Adam algorithm (Kingma & Ba, 2015) to adjust the learning rates. All methods involve with stochastic batches with singleton terms. For these terms, we use a batch size $B_1 = 64$. For the CVAEs, there are some pairwise terms of the sum involve sampling edges from the graph $G = (V, E)$ (i.e. the “positive sampling”) or sampling edges from the complete graph K_n (i.e. the “negative sampling”). We use a batch size $B_2 = 256$ for sampling these pairwise terms.