

Group Meeting - August 07, 2020

Paper review & Research progress

Truong Son Hy *

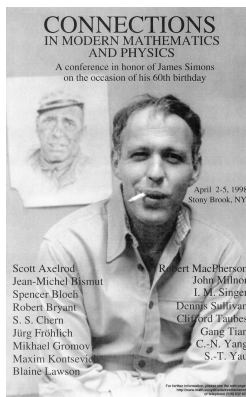
*Department of Computer Science
The University of Chicago

Ryerson Physical Lab



James Simons

I wasn't the fastest guy in the world. I wouldn't have done well in an Olympiad or a math contest. But I like to ponder. And pondering things, just sort of thinking about it and thinking about it, turns out to be a pretty good approach.



What I realize for myself

- ① I love Science and Research.
- ② I am no longer into Silicon valley, I don't want to work in the industry.
- ③ Reading scientific papers is more meaningful than working in a tech company.
- ④ Prof. Jordan Peterson: 'The pursuit of happiness is a pointless goal. We must instead search for meaning, not for its own sake, but as a defense against the suffering that is intrinsic to our existence.'



Correlated Variational Auto-Encoders, Da Tang, Dawen Liang, Tony Jebara, Nicholas Ruozzi (ICML 2019)



Standard VAEs (1)

Input data $\mathbf{x} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^D$. Standard VAEs assume that each data point \mathbf{x}_i is generated independently by the following process:

- 1 Generate the latent variables $\mathbf{z} = \{\mathbf{z}_1, \dots, \mathbf{z}_n\} \subseteq \mathbb{R}^d$ ($d \ll D$) by drawing i.i.d from the prior distribution (e.g. standard Gaussian distribution) $\mathbf{z}_i \sim p_0(\mathbf{z}_i)$, for each i .
- 2 Generate the data points $\mathbf{x}_i \sim p_\theta(\mathbf{x}_i | \mathbf{z}_i)$ from the model conditional distribution p_θ independently.



Standard VAEs (2)

Optimizing θ to maximize the likelihood p_θ requires computing **intractable** posterior distribution

$$p_\theta(\mathbf{z}|\mathbf{x}) = \prod_{i=1}^n p_\theta(\mathbf{z}_i|\mathbf{x}_i)$$

VAEs approximates this posterior distribution as $q_\lambda(\mathbf{z}|\mathbf{x}) = \prod_{i=1}^n q_\lambda(\mathbf{z}_i|\mathbf{x}_i)$ via amortized inference and maximize the evidence lower bound (ELBO):

$$\begin{aligned} L(\lambda, \theta) &= \mathbb{E}_{q_\lambda(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \mathcal{D}(q_\lambda(\mathbf{z}|\mathbf{x})||p_0(\mathbf{z})) \\ &= \sum_{i=1}^n \left[\mathbb{E}_{q_\lambda(\mathbf{z}_i|\mathbf{x}_i)}[\log p_\theta(\mathbf{x}_i|\mathbf{z}_i)] - \mathcal{D}(q_\lambda(\mathbf{z}_i|\mathbf{x}_i)||p_0(\mathbf{z}_i)) \right] \end{aligned}$$



Correlated priors on acyclic graphs (1)

Input data $\mathbf{x} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ with correlation structure given by an undirected graph $G = (V, E)$ in which $V = v_1, \dots, v_n$ is the set of vertices corresponding to each data point, and $(v_i, v_j) \in E$ if \mathbf{x}_i and \mathbf{x}_j are correlated. Define prior distribution p_0^{corr} of the latent variables $\mathbf{z}_1, \dots, \mathbf{z}_n$ over $(\mathbf{z}_1, \dots, \mathbf{z}_n) \in \mathbb{R}^d \times \dots \times \mathbb{R}^d$ whose singleton and pairwise marginal distributions satisfying:

$$p_0^{corr}(\mathbf{z}_i) = p_0(\mathbf{z}_i) \quad \forall v_i \in V$$

$$p_0^{corr}(\mathbf{z}_i, \mathbf{z}_j) = p_0(\mathbf{z}_i, \mathbf{z}_j) \quad \forall (v_i, v_j) \in E$$

Symmetry and marginalization consistency properties:

$$p_0(\mathbf{z}_i, \mathbf{z}_j) = p_0(\mathbf{z}_j, \mathbf{z}_i) \quad \forall \mathbf{z}_i, \mathbf{z}_j \in \mathbb{R}^d$$

$$\int p_0(\mathbf{z}_i, \mathbf{z}_j) d\mathbf{z}_j = p_0(\mathbf{z}_i) \quad \forall \mathbf{z}_i \in \mathbb{R}^d$$



Correlated priors on acyclic graphs (2)

The generative process of a CVAE:

- 1 Sample \mathbf{z} from the prior p_0^{corr} .
- 2 Sample each data point \mathbf{x}_i conditionally independently from \mathbf{z}_i .

Wainwright & Jordan, 2008:

$$p_0^{corr}(\mathbf{z}) = \prod_{i=1}^n p_0(\mathbf{z}_i) \prod_{(v_i, v_j) \in E} \frac{p_0(\mathbf{z}_i, \mathbf{z}_j)}{p_0(\mathbf{z}_i)p_0(\mathbf{z}_j)}$$

The prior is factorized as singletons and pairwise marginal distributions.



Singleton variational family: Approximate the posterior distribution $p(z|\mathbf{x})$ as

$$q_{\lambda}(z|\mathbf{x}) = \prod_{i=1}^n q_{\lambda}(z_i|\mathbf{x}_i)$$

Correlated variational family: $q_{\lambda}(z|\mathbf{x})$ is factorized as singleton and pairwise marginal distributions

$$q_{\lambda}(z|\mathbf{x}) = \prod_{i=1}^n q_{\lambda}(z_i|\mathbf{x}_i) \prod_{(v_i, v_j) \in E} \frac{q_{\lambda}(z_i, z_j|\mathbf{x}_i, \mathbf{x}_j)}{q_{\lambda}(z_i|\mathbf{x}_i)q_{\lambda}(z_j|\mathbf{x}_j)}$$



On general graphs (1)

Trivial generalization fails

For a general graph G :

$$p_0^{corr}(\mathbf{z}) = \prod_{i=1}^n p_0(\mathbf{z}_i) \prod_{(v_i, v_j) \in E} \frac{p_0(\mathbf{z}_i, \mathbf{z}_j)}{p_0(\mathbf{z}_i)p_0(\mathbf{z}_j)}$$

is not guaranteed to be a valid distribution.



On general graphs (1)

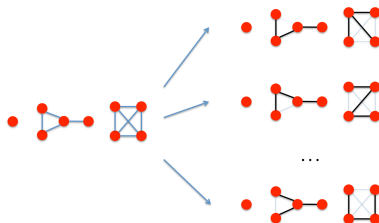
Trivial generalization fails

For a general graph G :

$$p_0^{corr}(\mathbf{z}) = \prod_{i=1}^n p_0(\mathbf{z}_i) \prod_{(v_i, v_j) \in E} \frac{p_0(\mathbf{z}_i, \mathbf{z}_j)}{p_0(\mathbf{z}_i)p_0(\mathbf{z}_j)}$$

is not guaranteed to be a valid distribution.

Solution: Factorized a general graph as a set of maximal acyclic graphs



On general graphs (2)

Maximal acyclic subgraph

For an undirected graph $G = (V, E)$, a subgraph $G' = (V', E')$ is a maximal acyclic subgraph of G if:

- 1 G' is acyclic.
- 2 $V' = V$, i.e., G' contains all vertices of G .
- 3 Adding any edge from E/E' to E' will create a cycle in G' .



On general graphs (2)

Maximal acyclic subgraph

For an undirected graph $G = (V, E)$, a subgraph $G' = (V', E')$ is a maximal acyclic subgraph of G if:

- 1 G' is acyclic.
 - 2 $V' = V$, i.e., G' contains all vertices of G .
 - 3 Adding any edge from $E \setminus E'$ to E' will create a cycle in G' .
- If G is connected, G' is a spanning tree of G .
 - Otherwise, a spanning forest.



On general graphs (3)

New prior distribution of \mathbf{z} as a uniform mixture over all subgraphs in \mathcal{A}_G :

$$p_0^{corr_g} = \frac{1}{|\mathcal{A}_G|} \sum_{G'=(V,E') \in \mathcal{A}_G} p_0^{G'}(\mathbf{z})$$

where

$$p_0^{G'}(\mathbf{z}) = \prod_{i=1}^n p_0(\mathbf{z}_i) \prod_{(v_i, v_j) \in E'} \frac{p_0(\mathbf{z}_i, \mathbf{z}_j)}{p_0(\mathbf{z}_i)p_0(\mathbf{z}_j)}$$

and \mathcal{A}_G is the set of all spanning forests of G .



On general graphs (4)

Log-likelihood:

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

where

$$q_{\lambda}^{G'}(\mathbf{z}|\mathbf{x}) = \prod_{i=1}^n q_{\lambda}(\mathbf{z}_i|\mathbf{x}_i) \prod_{(v_i, v_j) \in E'} \frac{q_{\lambda}(\mathbf{z}_i, \mathbf{z}_j|\mathbf{x}_i, \mathbf{x}_j)}{q_{\lambda}(\mathbf{z}_i|\mathbf{x}_i)q_{\lambda}(\mathbf{z}_j|\mathbf{x}_j)}$$



On general graphs (5)

- ① Singleton terms: all vertices have the same weight.
- ② Pairwise terms: an edge e 's weight is the fraction of times e appears among all subgraphs in \mathcal{A}_G .

$$w_{G,e}^{MAS} = \frac{|\{G' \in \mathcal{A}_G : e \in G'\}|}{|\mathcal{A}_G|}$$

Sum of all edge weights:

$$\sum_{e \in E} w_{G,e}^{MAS} = |V| - |CC(G)|$$

where $CC(G)$ is the set of connected components of G .

For a complete graph K_n , the weight $w_{K_n,e}^{MAS}$ for any edge e of K_n is $2/n$.



On general graphs (6)

New lower bound of the log-likelihood:

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



Some graph theory (1)

Matrix Tree Theorem (Chaiken & Kleitman, 1978)

For an undirected graph $G = (V, E)$, the number of spanning trees of G is the determinant of the sub-matrix of the Laplacian matrix L of G after deleting the i -th row and the i -th column, for any $i = 1, \dots, n$.

$$L_{i,j} = \begin{cases} \text{degree}(v_i) & \text{if } i = j, \\ -1 & \text{if } (v_i, v_j) \in E, \\ 0 & \text{otherwise.} \end{cases}$$



Some graph theory (2)

Number of spanning trees containing a particular edge

For an undirected graph $G = (V, E)$ and an edge $(v_i, v_j) \in E$, the number of spanning trees of G containing this edge is the determinant of the sub-matrix of the Laplacian matrix L of G after deleting the i -th, j -th rows and the i -th, j -th columns of it.

Complexity: $O(|E||V|^3)$ that is inefficient!



Some graph theory (3)

Smarter way

For an undirected connected graph $G = (V, E)$ and an edge $e = (v_i, v_j) \in E$, the weight $w_{G,e}^{MAS} = L_{i,i}^+ - L_{i,j}^+ - L_{j,i}^+ + L_{j,j}^+$ where L^+ is the **Moore-Penrose pseudo-inverse** of the Laplacian matrix L of G .

Complexity: $O(|V|^3)$



Algorithm to compute edge weights

Algorithm 1 Computing all weights $w_{G,e}^{\text{MAS}}$

Input: undirected graph $G = (V = \{v_1, \dots, v_n\}, E)$.
Compute all the connected components $\text{CC}_1, \dots, \text{CC}_K$
of G using depth-first search or breadth-first search.

for $k = 1$ **to** K **do**

 Compute the Moore-Penrose inverse L_k^+ of the Laplacian matrix L_k of the component CC_k .

 Apply Theorem 3 to compute $w_{G,e}^{\text{MAS}}$ for each edge e in the component CC_k .

end for

Return The weights $w_{G,e}^{\text{MAS}}$ for all $e \in E$.



Baselines and applications

Baselines:

- 1 Standard VAEs
- 2 GraphSAGE (Hamilton et al., 2017)

Experiments:

- 1 Bipartite correlation graph: MovieLens 20M dataset.
- 2 Perform spectral clustering on a synthetic dataset with a tree-structured latent variable graphical model.
- 3 Link prediction: Epinions dataset (Massa & Avesani, 2007).



Extension of this work: Markov Random Fields

Hammersley-Clifford theorem

A positive distribution $p(\mathbf{z}) > 0$ satisfies the CI (conditional independent) properties of an undirected graph G iff p can be represented as a product of factors, one per **maximal clique**, i.e.,

$$p(\mathbf{z}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{z}_c|\boldsymbol{\theta}_c)$$

where \mathcal{C} is the set of all the (maximal) cliques of G , and $Z(\boldsymbol{\theta})$ is the **partition function** given by:

$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{z}} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{z}_c|\boldsymbol{\theta}_c)$$

Note that the partition function is to ensure the overall distribution sums to 1.