

# Group Meeting - October 02, 2020

## Paper review & Research progress

Truong Son Hy \*

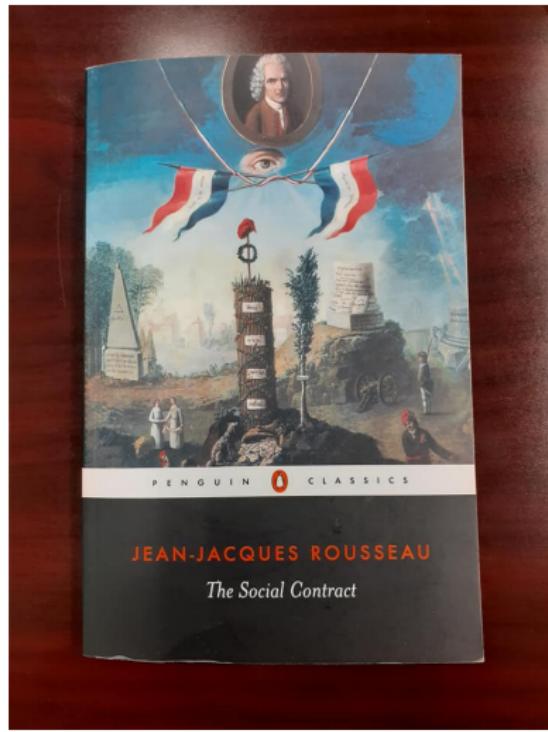
\*Department of Computer Science  
The University of Chicago

Ryerson Physical Lab



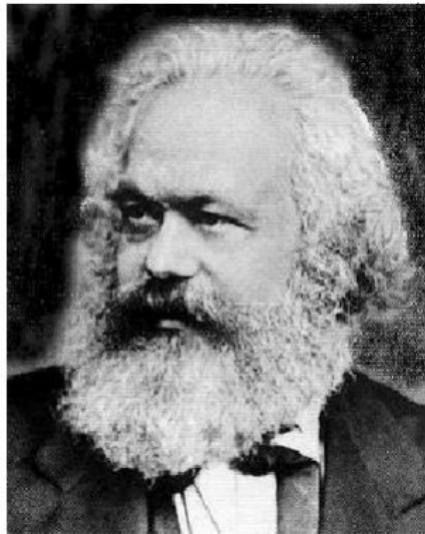
# Rousseau – *The Social Contract*

Man was born free, and he is everywhere in chains.



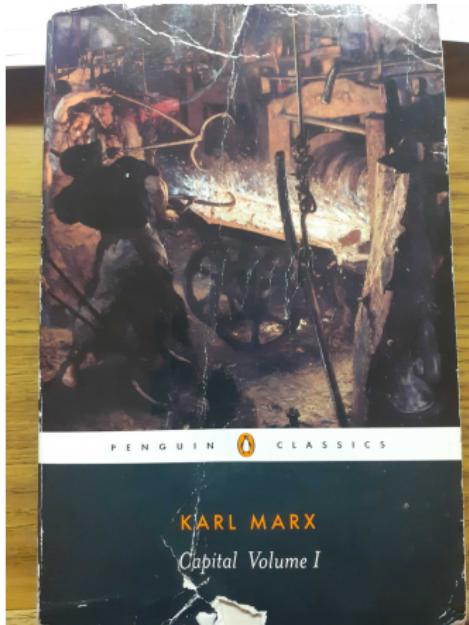
# Adam Smith – *The Wealth of Nations*

In general, if any branch of trade, or any division of labour, be advantageous to the public, the freer and more general the competition, it will always be the more so.



# Karl Marx – *Das Kapital*

As exchange-values, all commodities are merely definite quantities of congealed labour-time.



# Papers – Stochastic Variational Inference

- ① **Composing graphical models with neural networks for structured representations and fast inference** (NIPS 2016)  
<https://arxiv.org/pdf/1603.06277.pdf>
- ② **Structured VAEs: Composing Probabilistic Graphical Models and Variational Autoencoders** (spin-off from paper 1)  
[http://datta.hms.harvard.edu/wp-content/uploads/2018/01/pub\\_24.pdf](http://datta.hms.harvard.edu/wp-content/uploads/2018/01/pub_24.pdf)
- ③ **Neural Variational Inference for Text Processing** (ICML 2016)  
<https://arxiv.org/pdf/1511.06038.pdf>
- ④ **Learning to Generate with Memory** (ICML 2016)  
<https://arxiv.org/pdf/1602.07416.pdf>
- ⑤ **Generated Loss, Augmented Training, and Multiscale VAE**  
<https://arxiv.org/pdf/1904.10446.pdf>
- ⑥ **Learning Correlated Latent Representations with Adaptive Priors** (NeurIPS 2019 workshop)  
<http://bayesiandeeplearning.org/2019/papers/47.pdf>



# 3D Learning & Generation (planned for October 9)

- ① **PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation** (CVPR 2017)
- ② **Learning Representations and Generative Models for 3D Point Clouds** (ICML 2018)
- ③ **PointFlow: 3D Point Cloud Generation with Continuous Normalizing Flows** (ICCV 2019)
- ④ **3D Point Cloud Generative Adversarial Network Based on Tree Structured Graph Convolutions** (ICCV 2019)
- ⑤ **Learning Localized Generative Models for 3D Point Clouds via Graph Convolution** (ICLR 2019)
- ⑥ **Learning Gradient Fields for Shape Generation** (ECCV 2020)
- ⑦ **Progressive Point Cloud Deconvolution Generation Network** (ECCV 2020)
- ⑧ **Discrete Point Flow Networks for Efficient Point Cloud Generation** (ECCV 2020)



# 3D Learning & Generation (planned for October 9)

## Open question

What are we going to do with the 3D structure?

- For (molecular) Atom3D
- For generic (Computer Vision) point cloud

Partial (under construction) slides for October 09:

[http://people.cs.uchicago.edu/~hytruongson/  
Discussions-2020/Group\\_meeting\\_\\_\\_October\\_09\\_\\_2020.pdf](http://people.cs.uchicago.edu/~hytruongson/Discussions-2020/Group_meeting___October_09__2020.pdf)

Maybe I will present **Neural Turing Machines** to synchronize with the topic of **differentiable program(!)**



## Graph generation (as part of LibCCNs):

- Dynamic batching, NO patched/dummy vertices added
- Scale better
- Based on VAE with CCN encoder
- Tested with synthetic graphs (tree, star, ER)
- Tested with graph kernel benchmark datasets (DD, ENZYMES, MUTAG, PTC, etc.)
- [https://github.com/HyTruongSon/LibCCNs/tree/master/pytorch/small\\_graphs](https://github.com/HyTruongSon/LibCCNs/tree/master/pytorch/small_graphs)

## Room for CUDA extension:

- <https://github.com/HyTruongSon/LibCCNs/tree/master/cuda>



# Research – Theory Progress (1)

There are 2 ways to use the dipole decomposition:

Way 1:  
 $p(z|\theta) \propto \prod_{c \in \mathcal{C}} \psi_c(z_c|\theta_c) \leftarrow \text{Hanninen-Gilford theorem}$

$\mathcal{C}$  is the set of all dipoles

$$z_c = \{z_v | v \in \text{vec}\}$$

Procedure  
 $p(z|\theta) \propto \prod_{\text{over } c} \psi_c(z_{v_1}, z_v) \prod_v \psi_v(z_v)$

$$\psi_{ik}(z_{v_1}, z_v) = \exp\left(-\frac{1}{2} z_{v_1} \Delta_{ik} z_v^2\right)$$

$$\psi_k(z_v) = \exp\left(-\frac{1}{2} \Delta_{kk} z_v^2 + \beta_k z_v\right)$$

$$\Rightarrow p(z|\theta) \propto \exp\left[\sum_v -\frac{1}{2} z_v^T \Delta z_v\right]$$

Question: How to sample from higher-order  $p(z|\theta)$ .



Question: Statistical  
physics  $\rightarrow p^{(m)}$ .

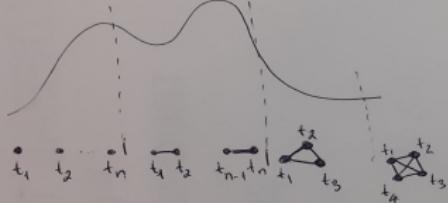
$$\psi_{ijk}(z_{v_1}, z_{v_2}, z_{v_3})$$

$$= \exp\left(\sum_{ijk} z_{v_1}^2 z_{v_2}^2 z_{v_3}^2\right)$$

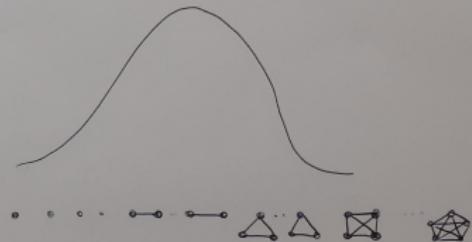
[Given  $\{z_v | v \in \text{vec}\}$ , the natural next  
decider takes care of the reconstruction.

## Research – Theory Progress (2)

Way 2: Distribution over all possible cliques.



vñ.



Question: Cliques overlapping in the construction.

# Papers 1

## **Composing graphical models with neural networks for structured representations and fast inference** (NIPS 2016)

Matthew J. Johnson, David Duvenaud, Alexander B. Wiltschko, Sandeep R. Datta, Ryan P. Adams

<https://arxiv.org/pdf/1603.06277.pdf>



# Composing graphical models – Proposals

## Proposals

- A general modeling and inference framework that combines:
  - ① **Probabilistic graphical models:** latent graphical models, message-passing algorithms.
  - ② **Deep learning methods:** neural network observation likelihoods, recognition networks to produce local evidence potentials.
- Single stochastic variational inference objective.
- Application to automatically segmenting and categorizing mouse behavior from raw depth video.



## Warped mixtures for arbitrary cluster shapes (1)

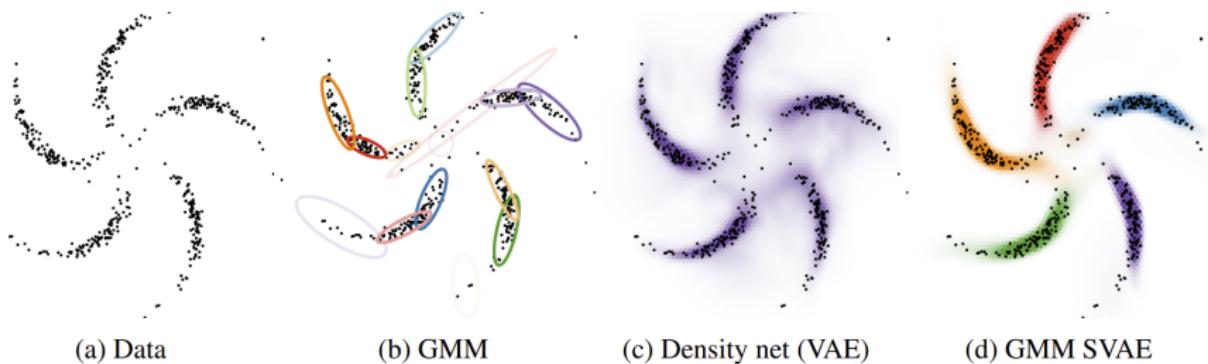


Figure 1: Comparison of generative models fit to spiral cluster data. See Section 2.1.

## Note:

- GMM is a generative model, but does not have the representation power of neural networks.
  - VAE does not have the power of a mixture model.



## Warped mixtures for arbitrary cluster shapes (2)

Consider the problem of modeling the data  $y = \{y_n\}_{n=1}^N$ . A standard approach to find the clusters in data is to fit a Gaussian Mixture Model (GMM):

$$\pi \sim \text{Dir}(\alpha)$$

$$(\mu_k, \Sigma_k) \sim_{\text{iid}} \text{NIW}(\lambda)$$

$$z_n | \pi \sim_{\text{iid}} \pi$$

$$y_n | z_n, \{(\mu_k, \Sigma_k)\}_{k=1}^K \sim_{\text{iid}} \mathcal{N}(\mu_{z_n}, \Sigma_{z_n})$$

### Note

For Dirichlet Process (DP), Dirichlet Process (Gaussian) Mixture Model (DPMM), and Latent Dirichlet Allocation (LDA), please check our past discussion in:

[http://people.cs.uchicago.edu/~hytruongson/  
Discussions-2020/Group\\_meeting\\_\\_\\_August\\_28\\_\\_2020.pdf](http://people.cs.uchicago.edu/~hytruongson/Discussions-2020/Group_meeting___August_28__2020.pdf)

# Warped mixtures for arbitrary cluster shapes (3)

Instead of using a GMM, a more flexible alternative would be a neural network density model:

$$\gamma \sim p(\gamma)$$

$$x_n \sim_{\text{iid}} \mathcal{N}(\mathbf{0}, \mathbf{1})$$

$$y_n | x_n, \gamma \sim_{\text{iid}} \mathcal{N}(\mu(x_n; \gamma), \Sigma(x_n; \gamma))$$

where  $\mu(x_n; \gamma)$  and  $\Sigma(x_n; \gamma)$  depend on  $x_n$  through some smooth (differential) parametric function (e.g. MLP), and  $p(\gamma)$  is a Gaussian prior.

## Limitation

The neural network density model fits the data well, but does not explicitly represent discrete mixture components.



# Warped mixtures for arbitrary cluster shapes (4)

Composing a latent GMM with nonlinear observations:

$$\pi \sim \text{Dir}(\alpha)$$

$$(\mu_k, \Sigma_k) \sim_{\text{iid}} \text{NIW}(\lambda)$$

$$\gamma \sim p(\gamma)$$

$$z_n | \pi \sim_{\text{iid}} \pi$$

$$x_n \sim_{\text{iid}} \mathcal{N}(\mu_{z_n}, \Sigma_{z_n})$$

$$y_n | x_n, \gamma \sim_{\text{iid}} \mathcal{N}(\mu(x_n; \gamma), \Sigma(x_n; \gamma))$$



# Latent (switching) linear dynamical systems (1)

## Note

For an introduction of linear dynamical systems including Kalman Filter and its variants (extended KF and unscented KF), please check the past discussion in:

[http://people.cs.uchicago.edu/~hytruongson/  
Discussions-2020/EKF\\_\\_\\_UJK.pdf](http://people.cs.uchicago.edu/~hytruongson/Discussions-2020/EKF___UJK.pdf)

## Sequence modeling

Consider the problem of **generatively modeling video** (that is a sequence of image frames):

- High-dimensional images  $\{y_n\}_{n=1}^N$ .
- Low-dimensional latent variables  $\{x_n\}_{n=1}^N$ .



# Latent (switching) linear dynamical systems (2)

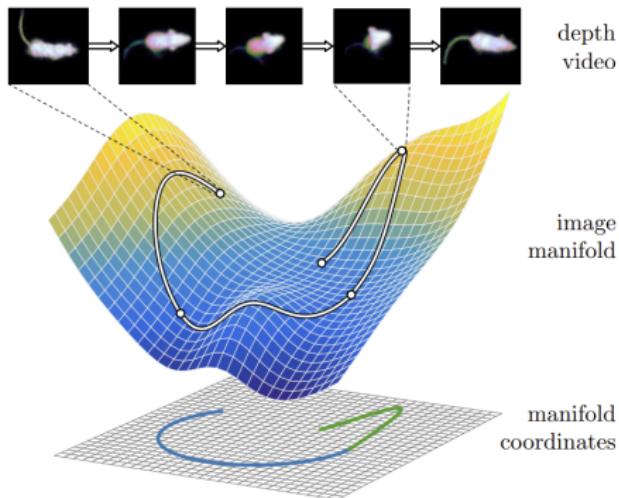


Figure 1. Illustration of unsupervised generative modeling applied to depth video of a mouse. Each video frame lies near a low-dimensional image manifold, and we wish to learn a parameterization of this manifold in which trajectories can be modeled with switching linear dynamics.



# Latent (switching) linear dynamical systems (3)

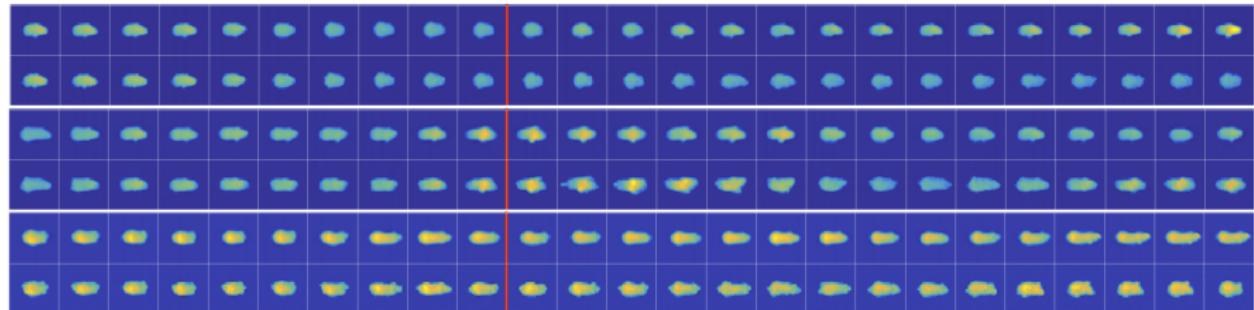


Figure 6: Predictions from an LDS SVAE fit to depth video. In each panel, the top is a sampled prediction and the bottom is real data. The model is conditioned on observations to the left of the line.



## Latent (switching) linear dynamical systems (4)

We introduce dependence through time between the latent Gaussian samples  $\{x_n\}_{n=1}^N$ . Each latent variable  $x_{n+1}$  depend on the previous latent variable  $x_n$  through a Gaussian linear dynamical system:

$$x_{n+1} = Ax_n + Bu_n, \quad u_n \sim_{\text{iid}} \mathcal{N}(\mathbf{0}, \mathbf{1}), \quad A, B \in \mathbb{R}^{m \times m}$$

We can construct a **latent switching linear dynamical system** (SLDS) to represent the data in terms of continuous latent states that evolve according to a discrete library of Markovian linear dynamics:

$$z_{n+1}|z_n, \pi \sim \pi_{z_n}, \quad x_{n+1} = A_{z_n}x_n + B_{z_n}u_n, \quad u_n \sim_{\text{iid}} \mathcal{N}(\mathbf{0}, \mathbf{1})$$

where  $n \in \{1, 2, \dots, N\}$  denotes time,  $z_n \in \{1, 2, \dots, K\}$  denotes discrete-valued latent state, and  $\pi = \{\pi_k\}_{k=1}^K$  denotes the Markov transition matrix and  $\pi_k \in \mathbb{R}_+^K$  is its  $k$ -th row.



# Latent (switching) linear dynamical systems (5)

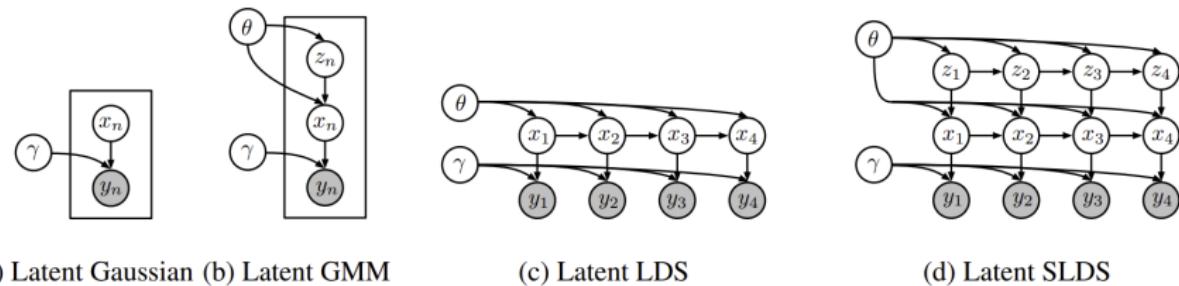


Figure 2: Generative graphical models discussed in Section 2.



# Structured mean field inference (1)

## Note

For an overview of variational inference in graphical models, please check our past discussion in:

[http://people.cs.uchicago.edu/~hytruongson/  
Discussions-2020/Group\\_meeting\\_\\_\\_August\\_14\\_\\_2020.pdf](http://people.cs.uchicago.edu/~hytruongson/Discussions-2020/Group_meeting___August_14__2020.pdf)

Consider learning a Gaussian linear dynamical system model with linear Gaussian observations:

$$x_n = Ax_{n-1} + Bu_{n-1}, \quad u_n \sim_{\text{iid}} \mathcal{N}(\mathbf{0}, \mathbf{1}), \quad y_n = Cx_n + Dv_n, \quad v_n \sim_{\text{iid}} \mathcal{N}(\mathbf{0}, \mathbf{1})$$

where:

- Latent states  $x = \{x_n\}_{n=1}^N$  and observations  $y = \{y_n\}_{n=1}^N$ .
- Parameters  $\theta = (A, B, C, D)$  with prior  $p(\theta)$ .



## Structured mean field inference (2)

To approximate the posterior  $p(\theta, x|y)$ , consider the mean field family  $q(\theta)q(x)$  and the variational inference objective:

$$\mathcal{L}[q(\theta)q(x)] = \mathbb{E}_{q(\theta)q(x)} \left[ \log \frac{p(\theta)p(x|\theta)p(y|x,\theta)}{q(\theta)q(x)} \right]$$

We can optimize the variational family  $q(\theta)q(x)$  to approximate the posterior  $p(\theta, x|y)$  by maximizing  $\mathcal{L}[q(\theta)q(x)]$ . For fixed  $q(\theta)$ ,

$$q^*(x) \triangleq \arg \max_{q(x)} \mathcal{L}[q(\theta)q(x)]$$

is computed by **exact** inference (message passing algorithms).

### Note

Regarding the factorization **restriction**  $p(\theta, x|y) \approx q(\theta)q(x)$ , please check Chapter 10 (Section 10.1.1) of Bishop's textbook.

# Structured variational autoencoders

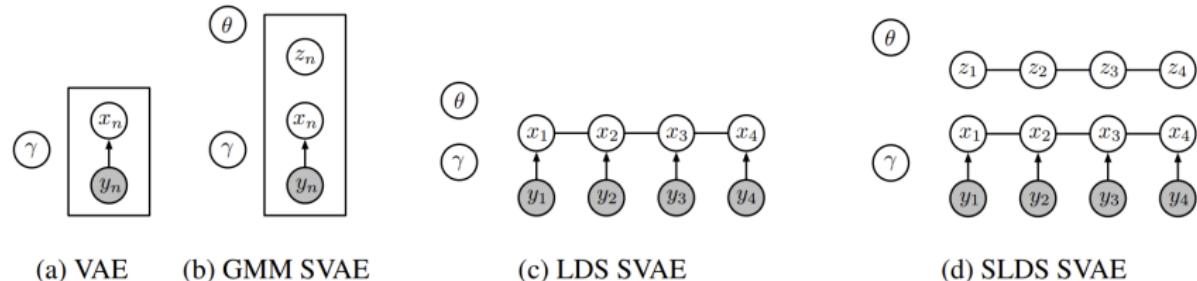


Figure 3: Variational families and recognition networks for the VAE [7] and three SVAE examples.

**Note:** Paper 1 & 2 overlap a lot, so we continue to discuss Structured VAEs with paper 2.

## More reference

[Matthew et al., 2014] **Stochastic Variational Inference for Bayesian Time Series Models** (ICML 2014)

<https://people.csail.mit.edu/mattjj/papers/icml2014svi.pdf>

## Structured VAEs: Composing Probabilistic Graphical Models and Variational Autoencoders (spin-off from paper 1)

Matthew J. Johnson, David Duvenaud, Alexander B. Wiltschko, Sandeep R. Datta, Ryan P. Adams

[http://datta.hms.harvard.edu/wp-content/uploads/2018/01/  
pub\\_24.pdf](http://datta.hms.harvard.edu/wp-content/uploads/2018/01/pub_24.pdf)

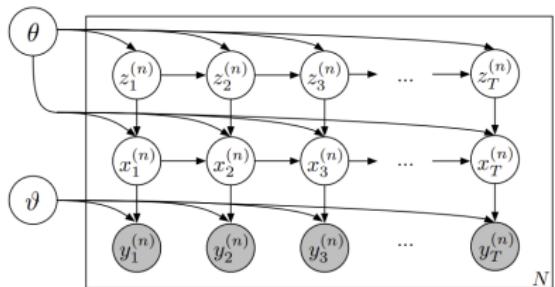


## Proposals

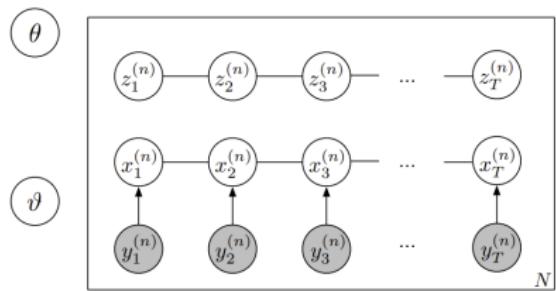
- A new framework for unsupervised learning that combines:
  - ➊ **Probabilistic graphical models:** to express structured probability distributions.
  - ➋ **Deep learning methods:** to learn flexible feature models.
- Scalable fitting algorithms for:
  - ➊ Natural gradient stochastic variational inference.
  - ➋ Graphical model message passing.
  - ➌ Backpropagation with the reparameterization trick.
- Application to structured time series model (e.g. mouse behavioral phenotyping).



# Conditional Random Field (1)



(a) SLDS generative model with nonlinear observation model parameterized by  $\vartheta$ .



(b) Structured CRF variational family with node potentials  $\{\psi(x_t^{(n)}; y_t^{(n)}, \phi)\}_{t=1}^T$  parameterized by  $\phi$ .

Figure 4. Graphical models for the SLDS generative model and corresponding structured CRF variational family.

Note: Between paper 1 & 2, there is a bit of notation change from  $\gamma$  to  $\vartheta$ .



## Conditional Random Field (2)

**[Wainwright & Jordan, 2008]** In mean field variational inference, one constructs a tractable variational family by **breaking dependencies** in the posterior. In this case, between the dynamics parameters  $\theta$ , the observation parameters  $\vartheta$ , the discrete state sequence  $z = z_{1:T}$  and the continuous state sequence  $x = x_{1:T}$ :

$$q(\theta, \vartheta, z_{1:T}, x_{1:T}) = q(\theta)q(\vartheta)q(z_{1:T})q(x_{1:T})$$

- **Naive** mean field model continues to break the dependencies among  $z_{1:T}$  and  $x_{1:T}$ .
- **Structured** mean field model **preserves** the correlations of these random variables **across time**.



# Conditional Random Field (3)

[Murphy, 2012] we parameterize the factor  $q(x_{1:T})$  as a Conditional Random Field (CRF) in terms of pairwise (edge) potentials and node potentials:

$$q(x_{1:T}) \propto \left( \prod_{t=1}^{T-1} \psi(x_t, x_{t+1}) \right) \left( \prod_{t=1}^T \psi(x_t; y_t, \phi) \right)$$

where  $\psi(x_t, x_{t+1})$  is the directed pairwise potential, and the node potential  $\psi(x_t; y_t, \phi)$  is a function of the observation  $y_t$ . **Remark:** without the direction, this is called as Markov Random Field (MRF).



# Conditional Random Field (4)

We choose each node potential to be a Gaussian factor that has (we use this alternative representation instead of mean vector and covariance matrix):

- Precision matrix  $J(y_t; \phi)$
- Potential vector  $h(y_t; \phi)$

depending on the corresponding observation  $y_t$  through a neural network (MLP):

$$\psi(x_t; y_t, \phi) \propto \exp \left\{ -\frac{1}{2} x_t^T J(y_t; \phi) x_t + h(y_t; \phi)^T x_t \right\}$$

$$[h(y_t; \phi), J(y_t; \phi)] \leftarrow \text{MLP}(y_t; \phi)$$

These **local recognition networks** allow us to fit a regression from each observation  $y_t$  to a probabilistic guess at the corresponding latent state  $x_t$ .



# Inference algorithm (1)

The optimization proceeds by alternating between optimizing  $q(x_{1:T})$  and optimizing  $q(z_{1:T})$ . First, fixing the factor on the discrete states  $q(z_{1:T})$ :

$$q^*(x_{1:T}) \propto \exp \left\{ \log \psi(x_1) + \sum_{t=1}^{T-1} \log \psi(x_t, x_{t+1}) + \sum_{t=1}^T \log \psi(x_t; y_t) \right\}$$

where:

$$\psi(x_1) = \mathbb{E}_{q(\theta)q(z)}[\log p(x_1|z_1, \theta)]$$

$$\psi(x_t, x_{t+1}) = \mathbb{E}_{q(\theta)q(z)}[\log p(x_{t+1}|x_t, z_t, \theta)]$$



## Inference algorithm (2)

Similarly, fixing  $q(x_{1:T})$  the optimal factor  $q^*(z_{1:T})$  is proportional to:

$$q^*(z_{1:T}) \propto \exp \left\{ \log \psi(z_1) + \sum_{t=1}^{T-1} \log \psi(z_t, z_{t+1}) + \sum_{t=1}^T \log \psi(z_t) \right\}$$

where:

$$\psi(z_1) = \mathbb{E}_{q(\theta)}[\log p(z_1|\theta)] + \mathbb{E}_{q(\theta)q(x)}[\log p(x_1|z_1, \theta)]$$

$$\psi(z_t, z_{t+1}) = \mathbb{E}_{q(\theta)}[\log p(z_{t+1}|z_t)]$$

$$\psi(z_t) = \mathbb{E}_{q(\theta)q(x)}[\log p(x_{t+1}|x_t, z_t, \theta)]$$



# Inference algorithm (3)

---

**Algorithm 2** Model inference subroutine for the SLDS

---

**Input:** Variational dynamics parameter  $\tilde{\eta}_\theta$  of  $q(\theta)$ , node potentials  $\{\psi(x_t; y_t)\}_{t=1}^T$  from recognition network

**function** INFERENCE( $\tilde{\eta}_\theta$ ,  $\{\psi(x_t; y_t)\}_{t=1}^T$ )

    Initialize factor  $q(x)$

**repeat**

$$q(z) \propto \exp\{\mathbb{E}_{q(\theta)q(x)} \ln p(z, x | \theta)\}$$

$$q(x) \propto \exp\{\mathbb{E}_{q(\theta)q(z)} \ln p(x | z, \theta)\} \prod_t \psi(x_t; y_t)$$

**until**  $q(z)$  and  $q(x)$  converge

$\hat{x} \leftarrow$  sample  $q(x)$

$$\bar{t}_{zx} \leftarrow \mathbb{E}_{q(z)q(x)} t_{zx}(z, x)$$

$$\text{KL} \leftarrow \text{KL}(q(\theta) \| p(\theta))$$

$$+ N \mathbb{E}_{q(\theta)} \text{KL}(q(z)q(x) \| p(z, x | \theta))$$

**return** sample  $\hat{x}$ , expected stats  $\bar{t}_{zx}$ , divergence KL

**end function**

---



## Inference algorithm (4)

To draw samples  $\hat{x}_{1:T} \sim q(x_{1:T})$ , we perform **message passing**. Given two neighboring nodes  $i$  and  $j$ , we define the message from  $i$  to  $j$ :

$$m_{i \rightarrow j}(x_j) \triangleq \int \psi(x_i; y_i) \psi(x_i, x_j) m_{k \rightarrow i}(x_i) dx_i$$

where  $k$  is the other neighbor of node  $i$ . Iteratively:

- ➊ Sample  $\hat{x}_1 \sim q(x_1)$  from:

$$q(x_1) \propto \psi(x_1) \psi(x_1; y_1) m_{2 \rightarrow 1}(x_1)$$

- ➋ Given sample  $\hat{x}_{t-1}$ , we sample  $\hat{x}_t \sim q(x_t | \hat{x}_{1:t-1})$  from:

$$q(x_t | \hat{x}_{1:t-1}) \propto \psi(\hat{x}_{t-1}, x_t) \psi(x_t; y_t) m_{t+1 \rightarrow t}(x_t)$$

- ➌ Finally, we sample  $\hat{x}_T \sim q(x_T | \hat{x}_{1:T-1})$  from:

$$q(x_T | \hat{x}_{1:T-1}) \propto \psi(\hat{x}_{T-1}, x_T) \psi(x_T; y_T)$$



# To our research: sequence modeling → graph modeling

## Hammersley-Clifford theorem

A positive distribution  $p(\mathbf{x}) > 0$  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.,

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

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

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

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

# Paper 3

**Neural Variational Inference for Text Processing (ICML 2016)**

Yishu Miao, Lei Yu, Phil Blunsom

<https://arxiv.org/pdf/1511.06038.pdf>



# Proposals

## Proposals

- A generic variational inference framework for **generative** and **conditional** models of text.
- The proposed neural variational document model combines a **continuous stochastic document representation** with a bag-of-words generative model. (**unsupervised document generation**)
- The proposed neural answer selection model employs a **stochastic representation layer within an attention mechanism** to extract the semantics between a question and answer pair. (**supervised question answering**)

## Latent Dirichlet Allocation (LDA)

Last time, we discussed about LDA (can be understood as 2-level Dirichlet Process Mixture Model) but it can only generate bag-of-words (without order).

[http://people.cs.uchicago.edu/~hytruongson/Discussions-2020/  
Group\\_meeting\\_\\_\\_August\\_28\\_\\_2020.pdf](http://people.cs.uchicago.edu/~hytruongson/Discussions-2020/Group_meeting___August_28__2020.pdf)

# Overview

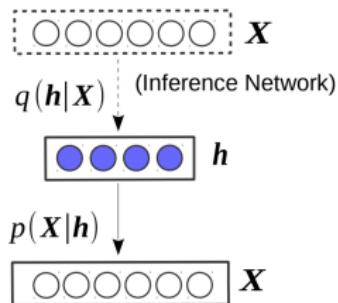


Figure 1. NVDM for document modelling.

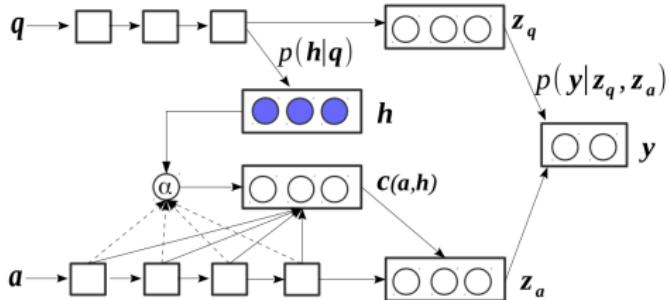


Figure 2. NASM for question answer selection.

NVDM = Neural Variational Document Model (**unsupervised**)  
NASM = Neural Answer Selection Model (**supervised**)

## Note

At first, I was looking for something like **Hierarchical VAEs** (e.g. normalizing flow), or **stochastic-in-time** (e.g. time-series) VAEs as paper 1 & 2. But this work is none of that. Furthermore, it can **only** generate documents as bags-of-words, no better than LDA.

# Neural Variational Inference Framework (1)

We introduce the common framework for both **unsupervised** NVDM and **supervised** NASM. Let  $\mathbf{h}$  denote the latent variable. Let  $\mathbf{x}$  and  $\mathbf{y}$  denote the observed parent of  $\mathbf{h}$ , and child nodes of  $\mathbf{h}$ , respectively.

$$p_{\theta}(\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{h}} p_{\theta}(\mathbf{y}|\mathbf{h})p_{\theta}(\mathbf{h}|\mathbf{x})p(\mathbf{x})$$

$$\log p_{\theta}(\mathbf{x}, \mathbf{y}) = \log \int \frac{q(\mathbf{h})}{q(\mathbf{h})} p_{\theta}(\mathbf{y}|\mathbf{h})p_{\theta}(\mathbf{h}|\mathbf{x})p(\mathbf{x})d\mathbf{h}$$

We can establish the lower bound:

$$\log p_{\theta}(\mathbf{x}, \mathbf{y}) \geq \mathbb{E}_{q(\mathbf{h})}[\log p_{\theta}(\mathbf{y}|\mathbf{h})p_{\theta}(\mathbf{h}|\mathbf{x})p(\mathbf{x}) - \log q(\mathbf{h})]$$



# Neural Variational Inference Framework (2)

Parameterized diagonal Gaussian  $\mathcal{N}(\mathbf{h}|\mu(\mathbf{x}, \mathbf{y}), \text{diag}(\sigma^2(\mathbf{x}, \mathbf{y})))$  as  $q_\phi(\mathbf{h}|\mathbf{x}, \mathbf{y})$ .  
The inference network is constructed as:

- ① Construct vector representations of the observed variables:  $\mathbf{u} = f_x(\mathbf{x})$ ,  $\mathbf{v} = f_y(\mathbf{y})$ .  $f_x$  and  $f_y$  are any type of neural networks.
- ② Assemble a joint representation:  $\pi = g(\mathbf{u}, \mathbf{v})$ .  $g$  is a concatenation and MLP.
- ③ Parameterize the variational distribution over the latent variable:  $\mu = l_1(\pi)$ ,  $\log \sigma = l_2(\pi)$ .  $l_1$  and  $l_2$  are linear transformations.
- ④ Reparameterization  $\mathbf{h} = \mu + \sigma \cdot \epsilon$ ,  $\epsilon^{(\ell)} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$ .



# Neural Variational Document Model

Let  $\mathbf{X} \in \mathbb{R}^{|V|}$  be the bag-of-words representation of a document, and  $\mathbf{x}_i \in \mathbb{R}^{|V|}$  be the one-hot representation of the  $i$ -th word. The lower bound:

$$\mathcal{L} = \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{X})} \left[ \sum_{i=1}^N \log p_\theta(\mathbf{x}_i|\mathbf{h}) \right] - \mathcal{D}_{\text{KL}}(q_\phi(\mathbf{h}|\mathbf{X}) || p(\mathbf{h}))$$

where  $N$  is the number of words in the document. The decoder  $p_\theta(\mathbf{x}_i|\mathbf{h})$  is modeled by multinomial logistic regression and shared across documents:

$$p_\theta(\mathbf{x}_i|\mathbf{h}) = \frac{\exp\{-E(\mathbf{x}_i; \mathbf{h}, \theta)\}}{\sum_{j=1}^{|V|} \exp\{-E(\mathbf{x}_j; \mathbf{h}, \theta)\}}$$

$$E(\mathbf{x}_i; \mathbf{h}, \theta) = -\mathbf{h}^T \mathbf{R} \mathbf{x}_i - \mathbf{b}_{x_i}$$

where  $\mathbf{R} \in \mathbb{R}^{K \times |V|}$  learns the semantic word embeddings and  $\mathbf{b}_{x_i}$  represents the bias term.



# Neural Answer Selection Model (1)

- **Problem setting:** A question  $q$  is associated with a set of answer sentences  $\{a_1, \dots, a_n\}$  together with their judgments  $\{y_1, \dots, y_n\}$ , where  $y_m = 1$  if the answer  $a_m$  is correct, and  $y_m = 0$  otherwise.
- **Task:** Predict  $y$  for the question-answer pair  $(q, a)$ .
- **Model:** In short, we employ 2 LSTMs (one for the question sentence, and the another one for the answer sentence) along with the attention mechanism.
- **Why don't we use a deterministic model instead?** The authors claimed this variational approach is more robust to noises(?!).



# Neural Answer Selection Model (2)

Variational lower bound:

$$\log p(\mathbf{y}|\mathbf{q}, \mathbf{a}) = \log \int p_{\theta}(\mathbf{y}|z_q(\mathbf{q}), z_a(\mathbf{a}, \mathbf{h}))p_{\theta}(\mathbf{h}|\mathbf{q})d\mathbf{h} \geq \\ \mathbb{E}_{q_{\phi}(\mathbf{h})}[\log p_{\theta}(\mathbf{y}|z_q(\mathbf{q}), z_a(\mathbf{a}, \mathbf{h}))] - \mathcal{D}_{\text{KL}}(q_{\phi}(\mathbf{h})||p_{\theta}(\mathbf{h}|\mathbf{q}))$$

where  $z_q(\mathbf{q})$  and  $z_a(\mathbf{a}, \mathbf{h})$  are the question and answer representations for predicting  $\mathbf{y}$ :

$$p_{\theta}(\mathbf{y} = 1|z_q, z_a) = \sigma(z_q^T \mathbf{M} z_a + b)$$

Note:  $z_q(\mathbf{q})$  is **deterministic**, while  $z_a(\mathbf{a}, \mathbf{h})$  is non-deterministic.



## Learning to Generate with Memory (ICML 2016)

Chongxuan Li, Jun Zhu, Bo Zhang

<https://arxiv.org/pdf/1602.07416.pdf>



# Proposals

## Proposals

- A **deep generative model** with a possibly large **external memory** and an **attention mechanism** to capture the local detail information (that is often lost in the bottom-up abstraction process).
- By adopting a **smooth** attention model, the whole network is trained end-to-end.
- Applications:
  - ① Density estimation
  - ② Image generation
  - ③ Missing value imputation



# Overall architecture

## Overall architecture

- Combination between the **stochastic layers** and **deterministic layers**.
- Each deterministic layer is associated with an **external memory** to capture local variant information.
- An attention mechanism is used to record information in the memory during learning, and retrieve information from the memory during data generation.

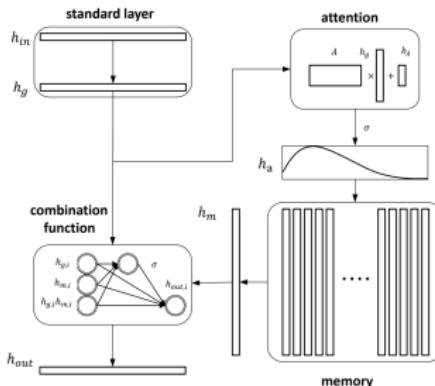


Figure 1. Architecture comparison between a standard layer (top-left part) and a layer with memory (the whole figure).



# Probabilistic DGMs with Memory (1)

## Note

This work is an extension to the **Hierarchical VAEs** to address its loss of detail information in the abstraction process, that we discussed in one of our meetings:

[http://people.cs.uchicago.edu/~hytruongson/  
Discussions-2020/Group\\_meeting\\_\\_\\_September\\_11\\_\\_2020.pdf](http://people.cs.uchicago.edu/~hytruongson/Discussions-2020/Group_meeting___September_11__2020.pdf)

Given a set of training data  $\mathcal{D}$ . The latent factors  $z_L, \dots, z_1$  are organized in a hierarchy:

- Draw the top-layer  $z_L \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$ .
- For  $\ell = L - 1, \dots, 0$ , calculate the mean parameters  $\mu_\ell = g_\ell(z_{\ell+1}; M_\ell)$  and draw the factors  $z_\ell \sim P_\ell(\mu_\ell)$ .
- $z_0 = x \leftarrow$  This is data.



## Probabilistic DGMs with Memory (2)

$$\mu_\ell = g_\ell(\mathbf{z}_{\ell+1}; M_\ell) \quad \mathbf{z}_\ell \sim P_\ell(\mu_\ell)$$

- $P_\ell$  is assumed to be of an exponential family form, with mean parameters  $\mu_\ell$ .
  - $g_\ell$  is a feed-forward deep neural network with  $I_\ell$  deterministic layers and a set of associated memories  $\{M_\ell^{(i)}\}_{i=0}^{I_\ell-1}$ , one per layer.
  - Each memory is a trainable matrix with dimension  $d_s \times n_s$ , where  $d_s$  is the number of slots in the memory, and  $n_s$  is the dimension of each slot.
  - The network  $g_\ell$  is parameterized as:
    - $\mathbf{h}_\ell^{I_\ell} = \mathbf{z}_{\ell+1} \leftarrow$  The top layer of  $g_\ell$ .
    - $\mathbf{h}_\ell^{(i)} = \phi(\mathbf{h}_\ell^{(i+1)}; M_\ell^{(i)}) \leftarrow$  The  $i$ -th layer of  $g_\ell$ .



# Inference (1)

Let  $\theta_g$  be the collection of parameters in the DGM.

$$p(\mathbf{x}, \mathbf{z}; \theta_g) = p(\mathbf{z}; \theta_g)p(\mathbf{x}|\mathbf{z}; \theta_g)$$

- **$Q$ -net (encoder):** To model the approximate posterior distribution  $q(\mathbf{z}|\mathbf{x}; \theta_r)$ , where  $\theta_r$  is the collection of parameters in the recognition model. **No external memory.**

$$\hat{\mathbf{h}}^{(i+1)} = \phi(V^{(i)}\hat{\mathbf{h}}^{(i)} + b_r^{(i)})$$

- **$P$ -net (decoder):** **With external memories.** The overall architecture is asymmetric.



# Inference (2)

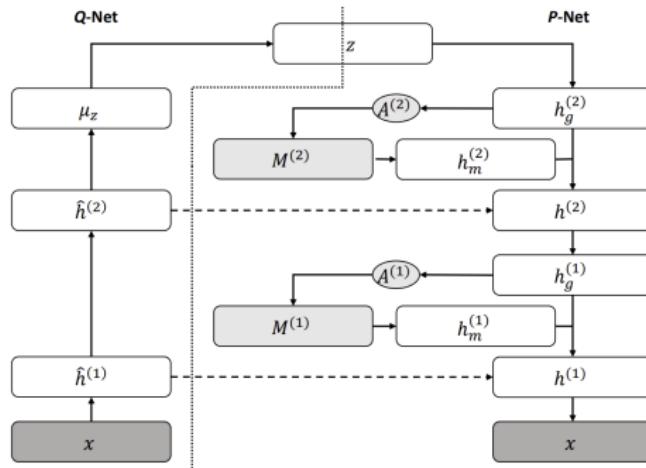


Figure 2. A model with one stochastic layer and two deterministic layers, where  $z$  is shared by the *P*-net and *Q*-net.

**Note:** The authors only used 1 stochastic layer, while the Hierarchical VAEs employs **several** stochastic layers. So they cannot claim it addresses the loss of information in Hierarchical VAEs.



## Inference (3)

A variational lower bound of log-likelihood for per data  $\mathbf{x}$  can be formulated as:

$$\mathcal{L}(\boldsymbol{\theta}_g, \boldsymbol{\theta}_r; \mathbf{x}) \triangleq \mathbb{E}_{q(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}_r)} [\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}_g) - \log q(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}_r)]$$

We add local reconstruction error terms as an optional regularizer, and jointly optimize the parameters in the generative model and the recognition model:

$$\min_{\boldsymbol{\theta}_g, \boldsymbol{\theta}_r} \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \left\{ \mathcal{L}(\boldsymbol{\theta}_g, \boldsymbol{\theta}_r; \mathbf{x}) + \sum_{i=1}^I \lambda^{(i)} \|\mathbf{h}^{(i)} - \hat{\mathbf{h}}^{(i)}\|_2^2 \right\}$$

where  $\lambda^{(i)}$  are the prefixed hyperparameters.



# How does this relate to our research?

Some idea we can use for **multiscale** graph generation:

- ① The recognition (encoder) model can apply **differential pooling** to construct a coarse-grained, bottom-up abstraction:
  - $\hat{\mathbf{h}}^{(1)}$  is a complete graph pooled from input graph  $\mathbf{x}$ :  $\hat{\mathbf{h}}^{(1)} = \text{POOL}(\mathbf{x})$
  - $\hat{\mathbf{h}}^{(i+1)} = \text{POOL}(\hat{\mathbf{h}}^{(i)})$
- ② The generative (decoder) model can reconstruct the graph in a reverse coarse-grained manner, from the high/general level into the low/detail level, along with a reconstruction loss:

$$\sum_{i=1}^I \lambda^{(i)} \|\mathbf{h}^{(i)} - \hat{\mathbf{h}}^{(i)}\|_2^2$$



## **Generated Loss, Augmented Training, and Multiscale VAE** **(only a technical report with VAE training tricks)**

Jason Chou, Gautam Hathi

<https://arxiv.org/pdf/1904.10446.pdf>



# Proposals

## Proposals

- **Discrete VAE** with tree recursive architecture shows limited capability of capturing field correlations within structured data.
- **Augmented training:** augmenting training data with generated variants.
- **Multiscale VAE:** training a VAE with multiple values of  $\beta$  simultaneously.



# Some insights

- **Mode collapse:** VAE is directly trained to encode all training examples and therefore is less prone to the failure mode of generating a **few memorized training examples** as in GAN.
- **Sajjadi et al., 2018:** VAE tends to have lower precision (e.g. blurry images for visual problems). **Rezende & Viola, 2018:** Holes in the latent space. (**This comes from the information bottleneck of VAE.**)
- **Higgins et al., 2017:** Fine-tuning the strength of the information bottleneck by the introduction of the KL-divergence weight  $\beta$ .
- **Manual annealing:**
  - ① KL-divergence warm-up
  - ② Controlled capacity increase



# Generalized $\beta$ -VAE

The evidence lower bound (ELBO):

$$\log p(x) \geq \mathbb{E}[\log p_\theta(x|z)] - \mathcal{D}_{\text{KL}}(q_\lambda(z|x)||p(z))$$

Modified ELBO:

$$\mathbb{E}[\log p_\theta(x|z)] - \beta \mathcal{D}_{\text{KL}}(q_\lambda(z|x)||p(z))$$

where  $\beta$  controls the strength of the information bottleneck on the latent vector. For higher values of  $\beta$ , we accept lossier reconstruction.

Higgins et al., 2017

This hyperparameter  $\beta$  is (theoretically justified) a Karush-Kuhn-Tucker (KKT) multiplier for maximizing  $\mathbb{E}[\log p_\theta(x|z)]$  under the inequality constraint  $\mathcal{D}_{\text{KL}}(q_\lambda(z|x)||p(z)) \leq c$  where  $c$  is a constant.



# Augmented training (1)

## Intuition

Augmented training extends the standard VAE training scheme. Instead of just taking one hop from the sampled latent vector and then minimizing the reconstruction loss, we take the second hop by taking the reconstructed sample as the new input to the encoder, and continue, and so on.

After `gen_start_step` steps:

1. Initialize  $n_{\text{augmented}}$  augmented latent vectors with sampled latent vectors of the current training batch.
2. Augment next training batch with variants generated from the augmented latent vectors.
3. After a training step, each augmented latent vector is replaced with either:
  - (a) The sampled latent vector of an example from the current training batch, selected without replacement, with probability  $p_{\text{sampled}}$ .
  - (b) The sampled latent vector of the variant generated from it with probability  $(1 - p_{\text{sampled}})$ .
4. Repeat from 2.



# Augmented training (2)

Mathematically saying:

$$x^{(1)} = g(z), \quad z \sim \mathcal{N}(\mu_\lambda(x), \sigma_\lambda(x))$$

$$x^{(2)} = g(z^{(1)}), \quad z^{(1)} \sim \mathcal{N}(\mu_\lambda(x^{(1)}), \sigma_\lambda(x^{(1)}))$$

...

$$x^{(n)} = g(z^{(n-1)}), \quad z^{(n-1)} \sim \mathcal{N}(\mu_\lambda(x^{(n-1)}), \sigma_\lambda(x^{(n-1)}))$$

The new objective:

$$\sum_{n=0}^{\infty} p_{\text{sampled}}^{\min(n,1)} (1 - p_{\text{sampled}})^{\max(n-1,0)} \left\{ \mathbb{E}[\log p_\theta(x^{(n)}|z^{(n)})] \right.$$

$$\left. - \beta \mathcal{D}_{\text{KL}}(q_\lambda(z^{(n)}|x^{(n)}) || p(z^{(n)})) \right\}$$



# Multiscale VAE (1)

## Observations:

- We can improve the model with KL-divergence warm-up (high  $\beta$ ) and cool-down (low  $\beta$ ).
- The objective functions with different  $\beta$  are **not** necessarily in conflict with each other.
- Training with  $\beta_{\max}$ , the highest possible  $\beta$  before  $q_\lambda(z|x)$  collapses into multivariate unit Gaussian, optimizes the global structure of the latent space.
- Training with smaller  $\beta$  optimizes the local structure.
- Training with  $\beta = 0$  optimizes autoencoding.



# Multiscale VAE (2)

New training scheme with multiple values of  $\beta$ , each trained on one worker:

1. Initialize  $n_{\text{KL\_weight}}$  standard deviation networks  $\sigma_{\lambda,i}(\mu)$ , where each standard deviation network is associated with a distinct but constant  $\beta$  value  $\beta_i$ . Without loss of generality, we assume  $\beta_i < \beta_{i+1}$ .
2. Assign  $(\beta_i, \sigma_{\lambda,i}(\mu))$  pairs evenly to workers, which otherwise share the same encoder / decoder.
3. Train the model with these workers.

In terms of objective function, we have

$$\sum_{i=0}^{\text{kl\_weight\_levels}-1} \mathbb{E}[\log p_\theta(x|z^{(i)})] - \beta_i KL(q_{\lambda,i}(z^{(i)}|x)||p(z^{(i)})) \quad (5)$$

where  $z^{(i)} \sim \mathcal{N}(\mu_\lambda(x), \sigma_{\lambda,i}(\mu_\lambda(x)))$

**Note:** This kind of **asynchronous** training is actually used in the industry to boost up the efficiency.



## **Learning Correlated Latent Representations with Adaptive Priors** (NeurIPS 2019 workshop)

Da Tang, Dawen Liang, Nicholas Ruozzi, Tony Jebara

<http://bayesiandeeplearning.org/2019/papers/47.pdf>



# Proposals

## Proposals

- Based on CVAEs, the authors proposed Adaptive Correlated Variational Auto-Encoders (ACVAEs): an **adaptive prior** distribution that can be **adjusted during training**.
- Applications: link prediction and hierarchical clustering.

## Correlated Variational Auto-Encoders (CVAEs)

For a presentation about the previous work (Tony Jebara's group) of this one, please check our previous discussion in:

[http://people.cs.uchicago.edu/~hytruongson/  
Discussions-2020/Group\\_meeting\\_\\_\\_August\\_07\\_\\_2020.pdf](http://people.cs.uchicago.edu/~hytruongson/Discussions-2020/Group_meeting___August_07__2020.pdf)



# 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  $p_{\theta}$  independently.



## Standard VAEs (2)

Optimizing  $\theta$  to maximize the likelihood  $p_\theta$  requires computing **intractable** posterior distribution

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

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

$$L(\lambda, \theta) = \mathbb{E}_{q_\lambda(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \mathcal{D}_{\text{KL}}(q_\lambda(\mathbf{z}|\mathbf{x})||p_0(\mathbf{z}))$$

$$= \sum_{i=1}^n \left[ \mathbb{E}_{q_\lambda(z_i|\mathbf{x}_i)}[\log p_\theta(\mathbf{x}_i|z_i)] - \mathcal{D}_{\text{KL}}(q_\lambda(z_i|\mathbf{x}_i)||p_0(z_i)) \right]$$



## Correlated prior (1)

Given an undirected correlation graph  $G = (V, E)$ , where  $(v_i, v_j) \in E$  represents that the data point  $x_i$  and  $x_j$  are correlated:

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

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

### Wainwright & Jordan, 2008

For acyclic graph, the prior is factorized as singletons and pairwise marginal distributions:

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

Our previous discussion about **[Wainwright & Jordan, 2008]**:

[http://people.cs.uchicago.edu/~hytruongson/Discussions-2020/  
Group\\_meeting\\_\\_\\_August\\_14\\_\\_2020.pdf](http://people.cs.uchicago.edu/~hytruongson/Discussions-2020/Group_meeting___August_14__2020.pdf)

## Correlated prior (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  $G'$ .

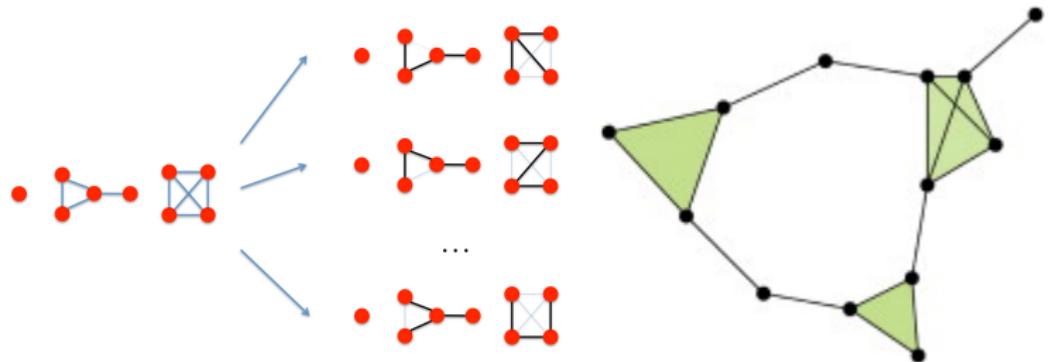
CVAEs optimize a different ELBO:

$$\log p_{\theta}(\mathbf{x}) = \mathbb{E}_{p_0^{corr}(\mathbf{z})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] = \frac{1}{|\mathcal{A}_G|} \sum_{G' \in \mathcal{A}_G} \mathbb{E}_{p_0^{G'}(\mathbf{z})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})]$$

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



# A different type of decomposition: maximal clique



The left is maximal acyclic graph (tree) decomposition.

The right is maximal clique decomposition.

**Note:** The right picture is taken from Wikipedia.



# Non-Uniform Mixture Prior (1)

Motivated by the ELBO:

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

Rather than using a uniform average, we employ a categorical distribution  $\pi \in \Delta^{|\mathcal{A}_G|-1}$  representing the normalized weights over all maximal acyclic subgraphs  $G' \in \mathcal{A}_G$  of  $G$ . With the non-uniform distribution  $\pi$ , we get a new ELBO:

$$\begin{aligned} \log p_{\pi, \theta}(\mathbf{x}) &\geq \mathbb{E}_{p_0^{\pi}(\mathbf{z})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] = \mathbb{E}_{G' \sim \pi} \left[ \mathbb{E}_{p_0^{\pi}(\mathbf{z})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] \right] \geq \\ &\geq \mathbb{E}_{G' \sim \pi} \left[ \mathbb{E}_{q_{\lambda}^{G'}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \mathcal{D}(q_{\lambda}^{G'}(\mathbf{z}|\mathbf{x}) || p_0^{G'}(\mathbf{z})) \right] \end{aligned}$$

with the non-uniform prior  $p_0^{\pi}(\mathbf{z}) = \mathbb{E}_{G' \sim \pi} [p_0^{G'}(\mathbf{z})]$ .



## Non-Uniform Mixture Prior (2)

If we optimize  $\pi$  together with all other parameters, we are optimizing with an adaptive prior:

$$\begin{aligned}\mathcal{L}^{\text{ACVAE}}(\pi, \lambda, \theta) := & \sum_{i=1} \left( \mathbb{E}_{q_\lambda(z_i|x_i)} [\log p_\theta(x_i|z_i)] - \text{KL}(q_\lambda(z_i|x_i)||p_0(z_i)) \right) - \sum_{(v_i, v_j) \in E} w_{G, \pi, (v_i, v_j)}^{\text{MAS}} \cdot \\ & \left( \text{KL}(q_\lambda(z_i, z_j|x_i, x_j)||p_0(z_i, z_j)) - \text{KL}(q_\lambda(z_i|x_i)||p_0(z_i)) - \text{KL}(q_\lambda(z_j|x_j)||p_0(z_j)) \right).\end{aligned}$$

where:

$$w_{G, \pi, e}^{\text{MAS}} = \sum_{G' \in \mathcal{A}_G, e \in E'} \pi(G')$$

ELBO with regularization:

$$\begin{aligned}\mathcal{L}^{\text{ACVAE-NS}}(\pi, \lambda, \theta) := & \mathcal{L}^{\text{ACVAE}}(\pi, \lambda, \theta) - \gamma \cdot \left( \sum_{i=1}^n \text{KL}(q_\lambda(z_i|x_i)||p_0(z_i)) \right. \\ & \left. + \frac{2}{n} \sum_{1 \leq i < j \leq n} \mathbb{E}_{q_\lambda(z_i, z_j|x_i, x_j)} \log \frac{q_\lambda(z_i, z_j|x_i, x_j)}{q_\lambda(z_i|x_i)q_\lambda(z_j|x_j)} \right).\end{aligned}$$



# Learning the Non-Uniform Mixture

Directly maximize  $\mathcal{L}^{\text{ACVAE}}(\pi, \lambda, \theta)$  with respect to  $\pi$ ,  $\lambda$  and  $\theta$ :

$$\max_{\lambda, \theta} \max_{\pi} \mathcal{L}^{\text{ACVAE}}(\pi, \lambda, \theta)$$

Alternatively, we can consider a minimax saddle-point optimization (which may lead to more robust inference):

$$\max_{\lambda, \theta} \min_{\pi} \mathcal{L}^{\text{ACVAE}}(\pi, \lambda, \theta)$$

## Proposition 1

If the optimization has global optima, then at least one optimum  $(\pi^*, \lambda^*, \theta^*)$  will have a  $\pi^*$  that places all of its probability mass on a single maximal acyclic subgraph  $G'^* \in \mathcal{A}_G$ .

That means at this optimum, the ELBO becomes the ELBO on a single spanning forest  $G'^*$ .

# ACVAEs Learning

---

**Algorithm 1** ACVAEs learning

---

**Input:** data  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^D$ , undirected graph  $G = (V = \{v_1, \dots, v_n\}, E)$ , parameter  $\gamma > 0$ .

Initialize the parameters  $\lambda, \theta$ . Initialize the weights  $w_{G, \pi, e}^{\text{MAS}}$  for each  $e \in E$ .

**while** not converged **do**

    Optimize the parameters  $(\lambda, \theta)$  using the gradients  $\nabla_{\lambda, \theta} \mathcal{L}^{\text{ACVAE}}(\pi, \lambda, \theta)$ .

    Compute the mass  $m_{(v_i, v_j)}$  for each edge  $e \in E$  with Eq. 11.

    Compute a minimum (if applying the empirical Bayes option in Eq. 9) or maximum (if applying the saddle-point option in Eq. 10) spanning tree of the graph with the masses  $m_{(v_i, v_j)}$ 's as the edge weights and update the weights  $w_{G, \pi, e}^{\text{MAS}}$  for each  $e \in E$  according to Eq. 12.

**end while**

**Return:** The parameters  $\lambda, \theta$ , the weights  $w_{G, \pi, e}^{\text{MAS}}$  for each  $e \in E$ .

---

The algorithm can be understood as:

**while** not converged **do**

- ① Update for  $\lambda$  and  $\theta$
- ② Update for  $\pi$

**end while**



# Metrics

Given learned latent embeddings  $\mathbf{z}_1, \dots, \mathbf{z}_n$ . We compute the pairwise distance  $d_{i,j}$  between each pair  $(\mathbf{z}_i, \mathbf{z}_j)$  as  $\|\mathbf{z}_i - \mathbf{z}_j\|_2^2$ . For the link prediction experiments, for each node  $v_i$ , we define the **Cumulative Reciprocal Rank** (CRR) as follows:

$$\text{CRR}_i = \sum_{(v_i, v_j) \in E_{\text{test}}} \frac{1}{|\{k : (v_i, v_k) \notin E_{\text{train}}, d_{i,k} \leq d_{i,j}\}|}$$

We further normalize CRR values to be in  $[0, 1]$  as normalized CRR (NCRR).  
**The larger CRR, the better.**

