



Variational Auto-Encoder (VAE)

2024. 10. 25.

MINGYU KIM



CONTENTS

1. Introduction
2. Preliminary
3. Auto-encoding variational bayes
4. Variational auto encoder
5. Experimental Result
6. What's Next

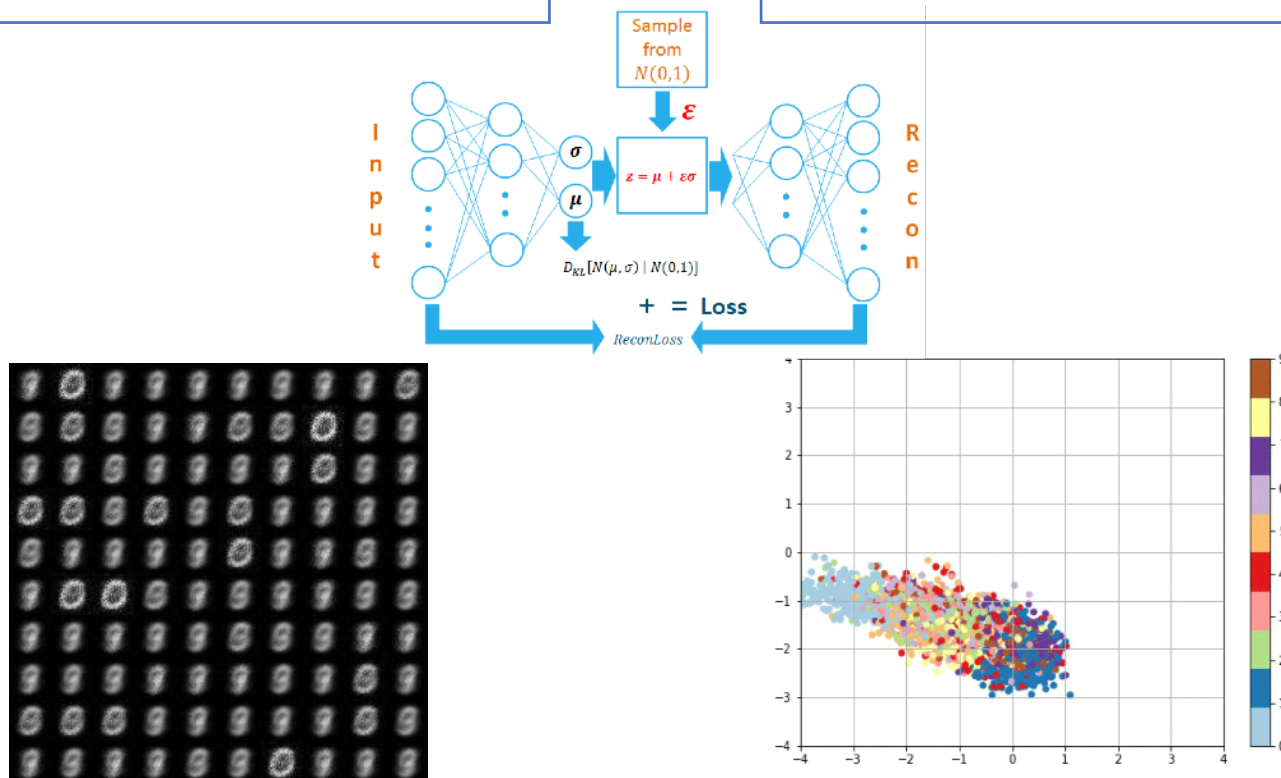
INTRODUCTION TO VAE (1)

Loss Function

- Variational Inference

Neural network architecture

- Auto encoder
- Mixture Density Network





INTRODUCTION TO VAE (2)

- **Characteristics**

- **Scalable** generative model
 - **Amortized** variational inference
 - Simple network architecture
 - Stochastic gradient descent (Back-propagation)
 - Scales to huge datasets
- **Inference** based on probability graphical model
 - Continuous Latent Space
 - Data reduction / Data Imputation
 - Liable to noisy



THIS SEMINAR

- **Paper “Auto-encoding variational bayes”**
 - “VAE” is a example of experiment results in this paper.
- **Introduce the rest of this paper without “VAE”**
 - Stochastic Gradient Variational Bayes(SGVB) estimator
 - Auto-encoding Variational Bayes algorithm
 - Reparameterization trick
 - Appendix
 - Full Variational Bayes
 - Marginal Likelihood estimator
- **Research trend of VAE**

PRELIMINARY

TRADITIONAL VARIATIONAL INFERENCE

BAYESIAN INFERENCE

- **General expression**

- A random variable $x \in X$
- A set $\{y_1, \dots, y_N\}$ of i.i.d observations

$$p(x, y_1, \dots, y_N) = p_0(x) \prod_{\{n=1\}}^N p(y_n | x)$$

- Given prior $p_{0(x)}$ and observations, posterior can be evaluated.

$$p(x | y_1, \dots, y_N) = \frac{p_0(x) \prod_{\{n=1\}}^N p(y_n | x)}{Z} := \frac{\bar{p}(x)}{Z}$$

$$Z := \int_X \bar{p}(x) dx$$

VARIATIONAL INFERENCE

- **Variational Inference**

- Approximate the target distribution $p(x)$ by using a tractable distribution $q(x)$ in Q

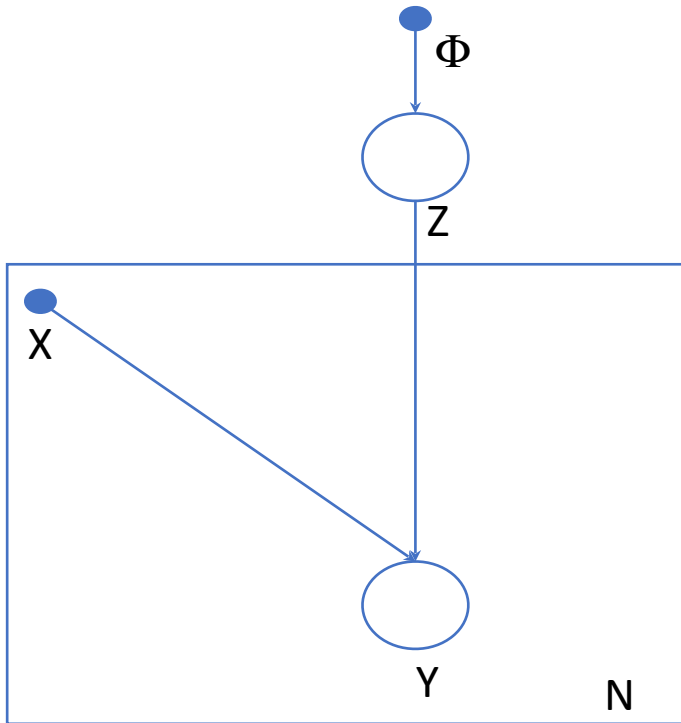
$$q^* = \operatorname{argmin}_{q \in Q} KL(q || p)$$

- **Variational Bound**

$$\begin{aligned} \log p(x) &= \int \log p(x) \cdot q(z) dz = \int \log \frac{p(x, z)}{p(z | x)} \cdot q(z) dz \\ &= \int \log \frac{p(x, z)}{q(z)} \cdot q(z) dz - \int \log \frac{q(z)}{p(z | x)} \cdot q(z) dz \end{aligned}$$

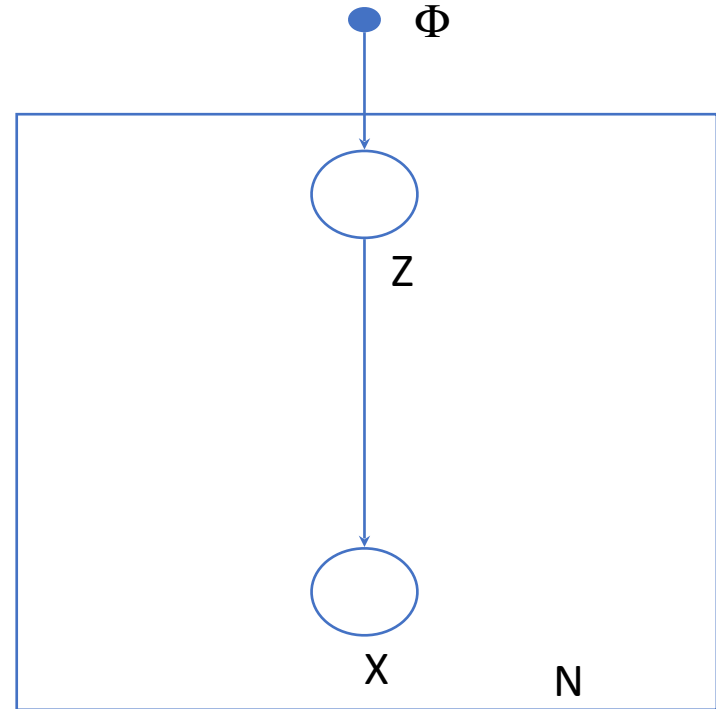
PROBABILITY GRAPHICAL MODEL

- **Supervised Learning**



e.g) Bayesian Logistic regression

- **Unsupervised Learning**
(Latent Variable Model)



e.g) EM algorithm for
Gaussian Mixture Model

PROBABILITY GRAPHICAL MODEL

- **Supervised learning**

e.g) Bayesian logistic regression

$$p(x | w, y) \cdot p(w) \propto p(w, x | y)$$

$$p(x | w, y) : \prod_{\{i=1\}}^n \sigma(w^T x)^{y_i} \cdot (1 - \sigma(w^T x))^{1-y_i} \rightarrow \text{Approximate Gaussian Dist.}$$

- **Latent Variable Model**

e.g) Gaussian Mixture Model (EM algorithm)

→ tractable conditional pdf : $p(z | x)$

$$q(z) \approx p(z | x, \theta^{old})$$

$$\theta^* = \operatorname{argmax}_{\theta} \sum_z p(z | x, \theta^{old}) \cdot \ln p(x, z | \theta)$$

TRADITIONAL VARIATIONAL INFERENCE

- **Mean field assumption**

- If $p(\mathbf{z} | \mathbf{x})$ is intractable ?

$$q(\mathbf{z}) = \prod_i q_i(z_i)$$

$$\ln q_j^*(z_j) = E_{i \neq j}[\ln p(\mathbf{x}, \mathbf{z})] + \text{const}$$

$$q_j^*(z_j) = \frac{\exp(E_{i \neq j}[\ln p(\mathbf{x}, \mathbf{z})])}{\int \exp(E_{i \neq j}[\ln p(\mathbf{x}, \mathbf{z})]) dz}$$

- Under specific probabilistic graphical model, approximate distribution can be evaluated by expectation of $E_{i \neq j}[\ln p(\mathbf{x}, \mathbf{z})]$
- Depending on each problem, we have to derive the form of $E_{i \neq j}[\ln p(\mathbf{x}, \mathbf{z})]$
- Through sequential optimization, All of $q_j(z_j)$ can be updated in order until converged.

TRADITIONAL VARIATIONAL INFERENCE (2)

- **VI with Stochastic gradient descent**

$$L = \int \log \frac{p(x, z)}{q_\phi(z)} \cdot q_\phi(z) dz$$

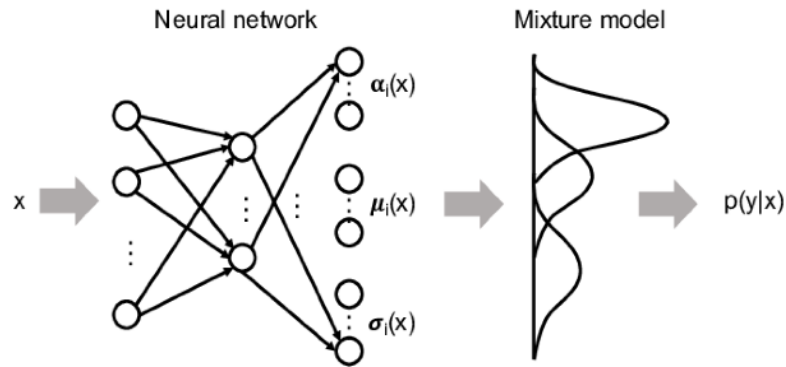
$$\nabla_\phi L = \int \ln p(x, z) \cdot \nabla_\phi \ln q(z) \cdot q(z) dz = E_{x \sim q}[\ln p(x, z) \cdot \nabla_\phi \ln q(z)]$$

$$\nabla_\phi L = E_{x \sim q}[\ln p(x, z) \cdot \nabla_\phi \ln q(z)] \approx \frac{1}{n} \sum_i \ln p(x_i, z) \cdot \nabla_\phi \ln q(z)$$

- In order to avoid exact derivation of expectation and sequential optimization, stochastic gradient descent method was introduced.
- This method only requires $\log p(x, z)$, $q_\phi(z)$ and its derivative
- All components can be easily handled
 - $\log p(x, z)$ can be derived from probabilistic graphical model
 - $q_\phi(z)$ and $\nabla_\phi q_\phi(z)$ are decided by users.
- However, since this estimator has a huge variance, it also requires many samples.

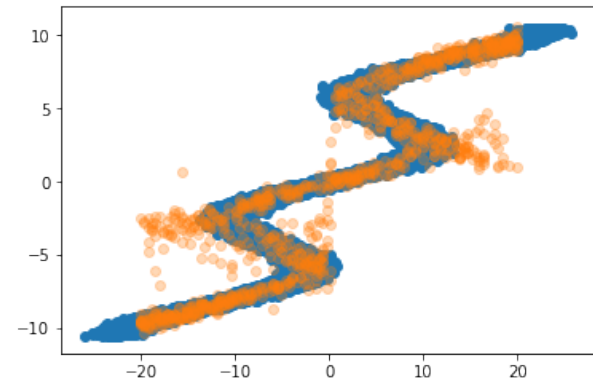
NETWORK FOR LEARNING PARAMETERS

- **Mixture density network**



[Structure of Mixture density Network]

$$E(w) = - \sum_{\{n=1\}}^N \ln \left\{ \sum_{\{k=1\}}^K \pi_k(x_n, w) N(t_n | u_k(x_n, w), \sigma_k^2(x_n, w)) \right\}$$



[Data point and result by Mixture Density Network]

Mixture density network presents that artificial neural network trains hyper-parameters for learning parameters of distributions



CHAPTER SUMMARY

- **Variational Inference**

- It provide the evidence in which mathematical analysis can be conducted
- Before emerging amortized VI, Most researches had a interest in efficient usage of VI using stochastic gradient descent and improving performance of approximate distribution for classification model
 - : Basic example : Bayesian logistic regression.
- Especially, VI with stochastic gradient descent influenced that “Auto-encoding variational bayes” were published.

- **Mixture Density Network**

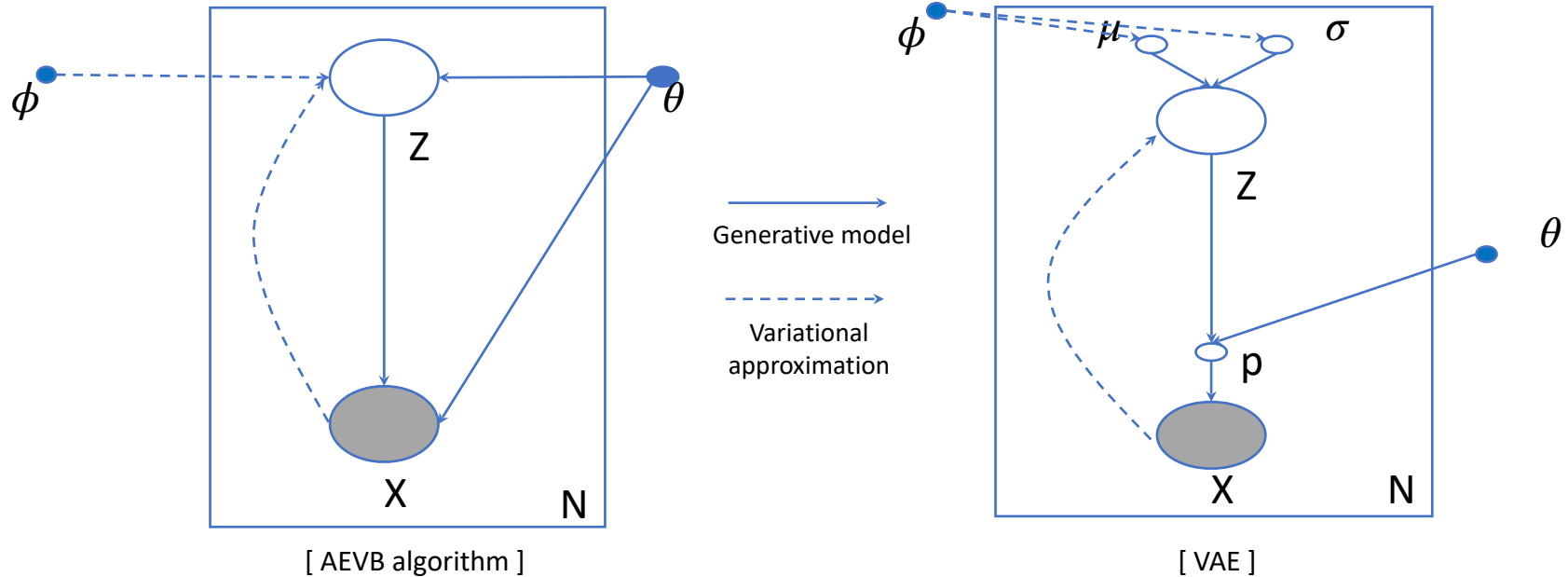
- Mixture Density Network is known to learn parameters of specific distributions using NN
- It is starting point to derive an approach to bayesian neural network

AUTO-ENCODING VARIATIONAL BAYES

PROBABILISTIC GRAPHICAL MODEL
STOCHASTIC VARIATIONAL GRADIENT BAYES ESTIMATOR

PROBABILISTIC GRAPHICAL MODEL

- Amortized generative model



- ϕ : hyper-parameters in recognition network
- θ : hyper-parameters in generative network

PROBLEMS SCENARIO

- **Traditional Latent variable model**

- Intractable : $p(z | x)$
- A large data set : Sampling based EM algorithm → Too slow

- **This paper has interested in**

- Using stochastic gradient descent
- Efficient approximate posterior inference of latent variable z

- **This paper introduce** (The reason why this method referred "Variational")

- $q_\phi(z | x)$: recognition model / approximate posterior distribution *to* $p_\theta(z | x)$
 - Figure out parameters of approximate distribution
- $p_\theta(x | z)$: generative model (defined as graphical model)
- Learning the recognition model parameter ϕ with generative model parameter θ using [reparameterization trick](#)

VARIATIONAL BOUND

- Developing variational Bound

$$\int \log \frac{p(x, z)}{q(z)} \cdot q(z) dz = \int \log \frac{p(x | z) \cdot p(z)}{q(z)} \cdot q(z) dz$$

$$\int (\log p(x | z) + \log p(z) - \log q(z)) \cdot q(z) dz$$

$$= \int \log p(x | z) \cdot q(z) dz + \int \log p(z) - \log q(z) dz$$

$$= \int \log p(x | z) \cdot q(z) dz - D_{KL}(q(z) || p(z))$$

STOCHASTIC GRADIENT VARIATIONAL BAYES ESTIMATOR (1)

- **Reparametrization Trick**

- Random variable $z \rightarrow$ Function of variables (Deterministic)
 - Capable of using neural network
 - $z \sim q(z|x) \rightarrow \bar{z} = g_\phi(x, \epsilon), \epsilon \sim p(\epsilon)$
 - Assume that $p(\epsilon)$ is defined ($p(\epsilon) \sim N(0,1)$)
 - The function g_θ is determined depending on $q(z|x)$
(In subsequent chapter, the forms of function g_θ is represented)
- Form of function g_θ
 - Tractable inverse CDF : similar to finding PDF
 - $F(x) = p, F^{-1}(p) = x$ (unique x) / $g_\phi(\epsilon, x) = F_\phi^{-1}(\epsilon; x) = \bar{Z}$
 - Location Scale Model :
 - $g_\phi(\epsilon, x) = \text{location} + \text{scale} \cdot \epsilon$

STOCHASTIC GRADIENT VARIATIONAL BAYES ESTIMATOR (2)

- **SGVB Estimator**

$$E_{q_{\phi}(z|x^i)}[f(z)] = E_{p(\epsilon)}\left[f\left(g_{\phi}(\epsilon, x^i)\right)\right] \approx \frac{1}{L} \sum_i f\left(g_{\phi}(\epsilon^l, x^i)\right) \text{ where, } \epsilon^l \sim p(\epsilon)$$

- $f(z) = \log p(x|z) + q(z) - p(z)$
 - Given x , ϵ sampled from $p(\epsilon)$, $z \approx \bar{z} = g_{\phi}(\epsilon, x^i)$, $f(z)$ can be evaluated.

- The form of $\frac{1}{L} \sum_i f\left(g_{\phi}(\epsilon^l, x^i)\right)$

- Stochastic gradient method can be introduced w.r.t ϕ, θ
(Back propagation)

STOCHASTIC GRADIENT VARIATIONAL BAYES ESTIMATOR (3)

- **SGVB Estimator A**

$$\tilde{L}^A(\theta, \phi; x^i) = \frac{1}{L} \sum_l \log p_\theta(x^i, z^{i,l}) - \log q_\phi(z^{i,l} | x^i) \text{ where, } g_\phi(\epsilon, x^i), \epsilon^l \sim p(\epsilon)$$

- $\log p_\theta(x^i, z^{i,l})$ is defined by probabilistic graphical model
- $\log q_\phi(z^{i,l} | x^i)$ can be determined $g_\phi(\epsilon, x^i), \epsilon^l \sim p(\epsilon)$

- **SGVB Estimator B**

$$\tilde{L}^B(\theta, \phi; x^i) = \frac{1}{L} \sum_l \log p_\theta(x^i | z^{i,l}) - D_{KL}(q_\phi(z | x^i) || p_\theta(z))$$

where, $g_\phi(\epsilon, x^i), \epsilon^l \sim p(\epsilon)$

- $D_{KL}(q_\phi(z | x^i) || p_\theta(z))$ is analytically evaluated (Gaussian Dist.)
- $D_{KL}(q_\phi(z | x^i) || p_\theta(z))$ doesn't require sample z , only require parameter of approximate dist.

AUTO ENCODING VARIATIONAL BAYES ALGORITHM(1)

Algorithm 1 Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings $M = 100$ and $L = 1$ in experiments.

$\theta, \phi \leftarrow$ Initialize parameters

repeat

$\mathbf{X}^M \leftarrow$ Random minibatch of M datapoints (drawn from full dataset)

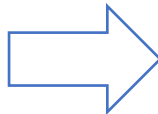
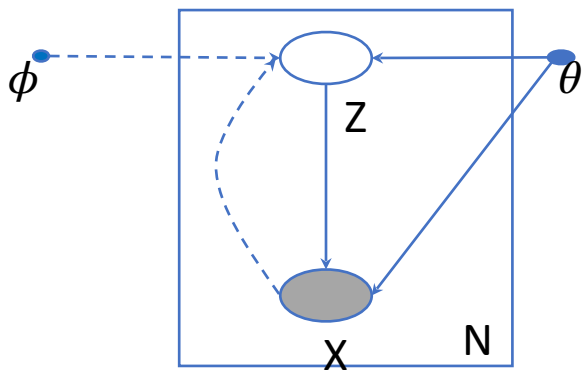
$\epsilon \leftarrow$ Random samples from noise distribution $p(\epsilon)$

$\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$ (Gradients of minibatch estimator (8))

$\theta, \phi \leftarrow$ Update parameters using gradients \mathbf{g} (e.g. SGD or Adagrad [DHS10])

until convergence of parameters (θ, ϕ)

return θ, ϕ



$$\ln p_{\theta}(X, Z) = \sum_i \ln p_{\theta}(x^i | z^i)$$

$$\ln q_{\phi}(Z) = \sum_i \ln q_{\phi}(z^i)$$

$$\ln p_{\theta}(Z) = \sum_i \ln p_{\theta}(z^i)$$

AUTO ENCODING VARIATIONAL BAYES ALGORITHM (2)

- **AEVB for SGVB estimator**

$$\frac{1}{M} \sum_i \tilde{L}^A(\theta, \phi; x^i) = \frac{1}{M} \sum_i \frac{1}{L} \sum_l \log p_{\theta}(x^i, z^{i,l}) - \log q_{\phi}(z^{i,l} | x^i)$$

$$\frac{1}{M} \sum_i \tilde{L}^B(\theta, \phi; x^i) = \frac{1}{M} \sum_i \frac{1}{L} \sum_l \log p_{\theta}(x^i | z^{i,l}) - D_{KL}(q_{\phi}(z | x^i) || p_{\theta}(z))$$



CHAPTER SUMMARY

- **SGVB Estimator**

- This estimator was developed to use an approach of neural network for learning probabilistic latent variable model
- Reparameterization trick is the key point to derive SGVB estimator

- **AEVB algorithm**

- It can train a neural network using gradient of loss function (SGVB estimator) w.r.t hyper-parameters in recognition and generative network at once.
- This gradient can be evaluated by back-propagation

VARIATIONAL AUTO ENCODER

FRAMEWORK

NETWORK ARCHITECTURE

FULL VARIATIONAL BAYES

MARGINAL LIKELIHOOD ESTIMATOR

FRAMEWORK (1)

- **Loss function** (Bernoulli Case)

- $p(z) \sim N(0, I)$
- $p_\theta(x | z) \sim N(u_\theta(z), \sigma_\theta(z))$ or $Bern(p_\theta(z))$
- $q_\phi(z | x) \sim N(u_\phi(x), \text{diag}_\phi(x))$

- **Estimator (Loss function)**

$$\tilde{L}^B(\theta, \phi; x^i) = \frac{1}{L} \sum_l \log p_\theta(x^i | z^{i,l}) - D_{KL}(q_\phi(z | x^i) || p_\theta(z))$$

where $g_\phi(\epsilon, x^i) = u_\theta(x) + \epsilon \cdot \sigma_\theta(x)$

$\epsilon \sim N(0, I), \quad l = 1$

- Since $l = 1$, neural network architecture can be similar to “auto-encoder”

FRAMEWORK (2)

- **Loss function**

- $p_{\theta}(x|z) \sim \text{Bern}(p_{\theta}(z))$

$$\rightarrow \log p_{\theta}(x|z) \sim \sum_i x_i \log p_{\theta}(z_i) + (1 - x_i) \log(1 - p_{\theta}(z_i))$$

- $D_{KL}(q_{\phi}(z|x^i) || p_{\theta}(z))$

- $p(z) \sim N(0, I), q_{\phi}(z|x) \sim N(u_{\phi}(x), \text{diag}_{\phi}(x))$

- $D_{KL}(q_{\phi}(z|x) || p_{\theta}(z)) =$

$$\frac{1}{2} \{ (0 - u_{\phi}(x))^T I^{-1} (0 - u_{\phi}(x)) + \text{trace}(I^{-1} \text{diag}_{\phi}(x)) - k + \ln \frac{|I|}{|\text{diag}_{\phi}(x)|} \}$$

$$= \sum_j u_{\phi}(x)_{(j)}^2 - 1 + \sigma_{\phi}^2(x)_{(j)} - \ln \sigma_{\phi}^2(x)_{(j)}$$

where $j = \text{index of dimension of latent variable } z$

FRAMEWORK (3)

- **Loss Function**

$$\therefore \sum_i x_i \log p_{\theta}(z_i) + (1 - x_i) \log(1 - p_{\theta}(z_i)) +$$

$$\sum_i \sum_j u_{\phi}(x_i)_{(j)}^2 - 1 + \sigma_{\phi}^2(x_i)_{(j)} - \ln \sigma_{\phi}^2(x_i)_{(j)}$$

STRUCTURE OF VAE(1)

- **Network Architecture**

- Recognition network (like density network)

$$u_{\phi}(x_i) = w_2 \tanh(w_1 x_i + b_1) + b_2$$

$$\sigma_{\phi}(x_i) = w_4 \tanh(w_3 x_i + b_3) + b_4$$

- Generative network

- Bernoulli Case

$$p_{\theta}(z_i) = f_a(w_2 \tanh(w_1 z_i + b_1) + b_2)$$

where $f_a(\cdot)$: element – wise sigmoid function

$$\theta = \{w_1, b_1, w_2, b_2\}$$

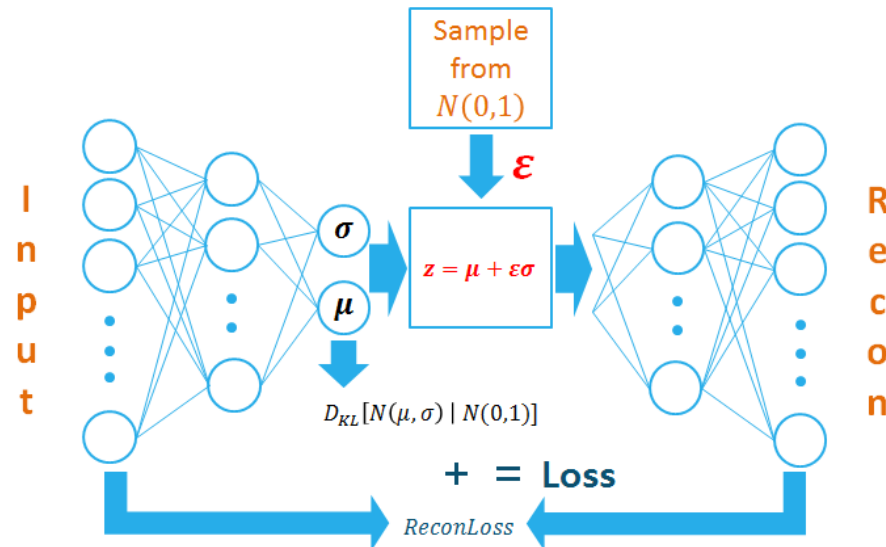
- Gaussian Case

$$u_{\theta}(z_i) = w_4 (\tanh(w_3 z + b_3) + b_4)$$

$$\sigma_{\theta}(z_i) = w_5 (\tanh(w_3 z + b_3) + b_5)$$

$$\theta = \{w_3, w_4, w_5, b_3, b_4, b_5\}$$

VAE



- **Probabilistic modeling**

- All data should be mapped to latent space defined by probability graph model
- Robust against some noises
 - Key difference between VAE and general AE

MARGINAL LIKELIHOOD ESTIMATOR (1)

- **Derivation**

- Marginal Likelihood : $p_{\theta}(x)$
- Marginal likelihood estimator

$$\begin{aligned}\frac{1}{p_{\theta}(x_i)} &= \int \frac{q(z)}{p_{\theta}(x_i)} dz = \frac{\int q(z) \cdot \frac{p_{\theta}(x_i, z)}{p_{\theta}(x_i, z)} dz}{p_{\theta}(x_i)} = \int \frac{p_{\theta}(x_i, z)}{p_{\theta}(x_i)} \cdot \frac{q(z)}{p_{\theta}(x_i, z)} dz \\ &= \int p_{\theta}(z | x_i) \frac{q(z)}{p_{\theta}(x_i, z)} dz \approx \frac{1}{L} \sum_l \frac{q_{\phi}(z^{(l)})}{p_{\theta}(z) p_{\theta}(x_i | z^{(l)})}\end{aligned}$$

- In VAE, marginal likelihood estimator cannot be evaluated because it requires $l > 1$
- If it was forced to evaluate the marginal likelihood in VAE, the result is almost same as loss function, it doesn't give any special meaning.

MARGINAL LIKELIHOOD ESTIMATOR (2)

- **Gradient MCMC**

- Max Welling and published "Bayesian learning via stochastic gradient langevin dynamics (2011), ICML2011"

$$\Delta\theta_t = \frac{\epsilon_t}{2}(\nabla \log p(\theta_t) + \frac{N}{n} \sum_{i=1}^n \nabla \log p(x_i | \theta_t)) + \eta_t$$

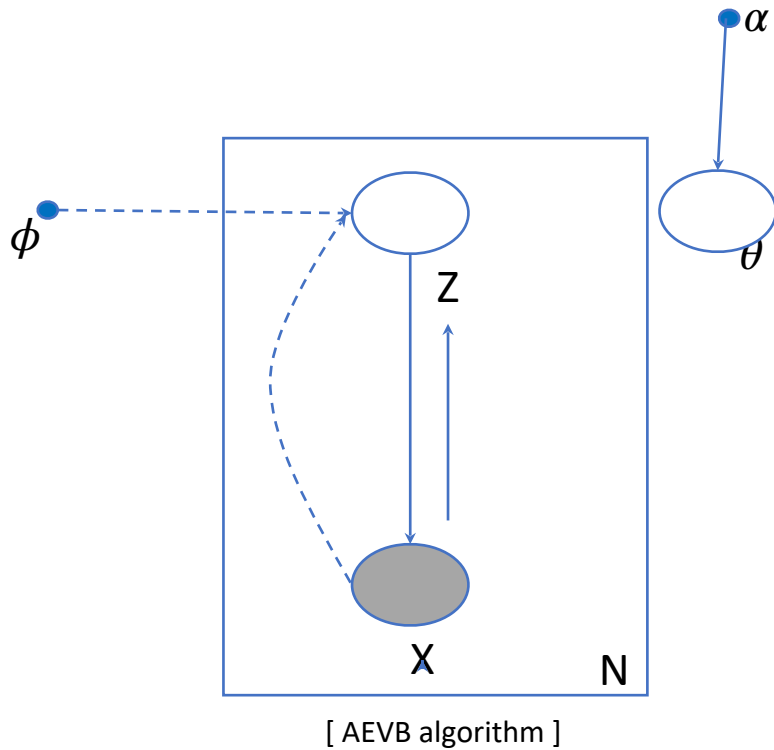
where $\eta_t \sim N(0, \epsilon_t)$

Update (stochastic gradient descent)

$$\theta = \theta_t + \epsilon'_t \cdot \Delta\theta_t$$

- In this paper, $\nabla \log p(z_t)$, $\nabla \log p(x_i | z_t)$ can be evaluated.
- Under this method, marginal likelihood estimator is established.

FULL VARIATIONAL BAYES



- **Framework**

- $p_{\theta}(z) \sim N(0, I)$, where $z \in R^J$, $I \in R^{J \times J}$
- $p_{\alpha}(\theta) \sim N(0, I)$, where $\theta \in R^M$, $I \in R^{M \times M}$
- $p_{\theta}(x | z) \sim N(u_{\theta}(z), \sigma_{\theta}(z))$ or $Bern(p_{\theta}(z))$
- $q_{\phi}(z | x) \sim N(u_{\phi}(x), \text{diag}_{\phi}(x))$

FULL VARIATIONAL BAYES

- **Variational Bound**

$$\begin{aligned} \ln p_{\theta}(X) &\approx \int \left(\ln P_{\alpha}(X | \theta) + \ln P_{\alpha}(\theta) - \ln q_{\phi}(\theta) \right) \cdot q_{\theta}(\theta) d\theta \\ &= \frac{1}{L} \sum_{l=1}^L E \left[\ln P_{\alpha}(X | \theta) \right] + \frac{1}{2} \sum_{M=1}^M \{ 1 + \ln \sigma_{\theta_m}^{2(l)} - \mu_{\theta_m}^{2(l)} - \sigma_{\theta_m}^{2(l)} \} \end{aligned}$$

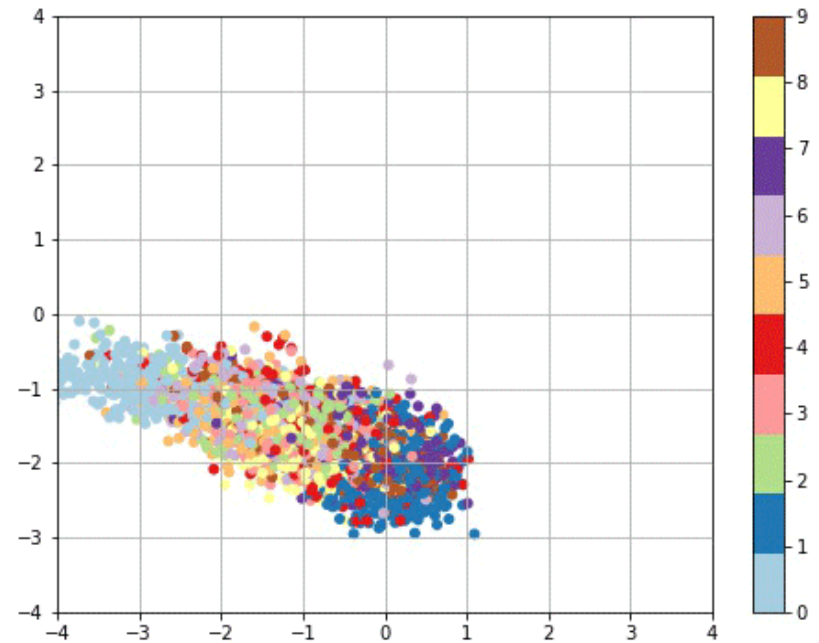
$$\begin{aligned} \ln P_{\alpha}(X | \theta) &= \int \log \frac{p(x, z)}{q(z)} \cdot q(z) dz - \int \log \frac{q(z)}{p(z | x)} \cdot q(z) dz \\ &\approx \frac{1}{K} \sum_{i=1}^N \left(\sum_{k=1}^K \ln p_{\theta}(x_i | z_{i,j}) \right) + \frac{1}{2} \sum_{j=1}^J \{ 1 + \ln \sigma_{z_{i,j}}^{2(k)} - \mu_{z_{i,j}}^{2(k)} - \sigma_{z_{i,j}}^{2(k)} \} \end{aligned}$$

- In this paper, $L = K = 1$ and $\ln p_{\theta}(x_i | z_{i,j}) = \ln p_{\theta}(x | z), \forall x_i$

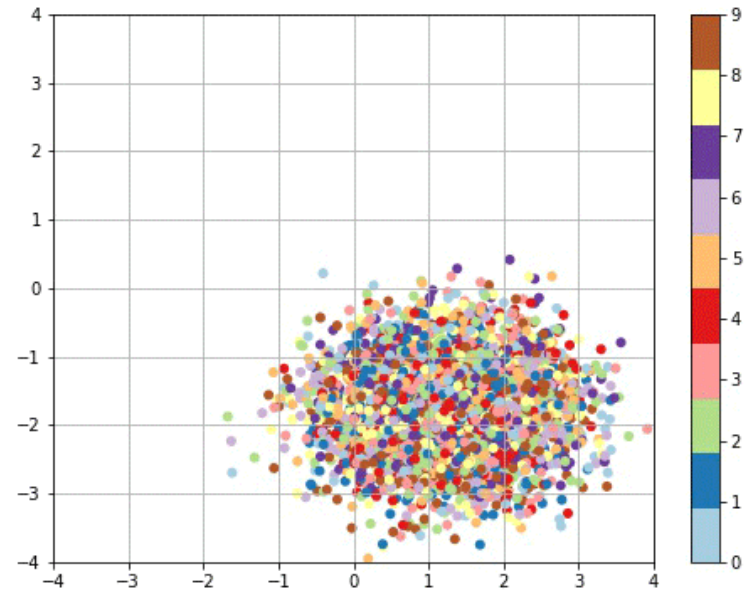
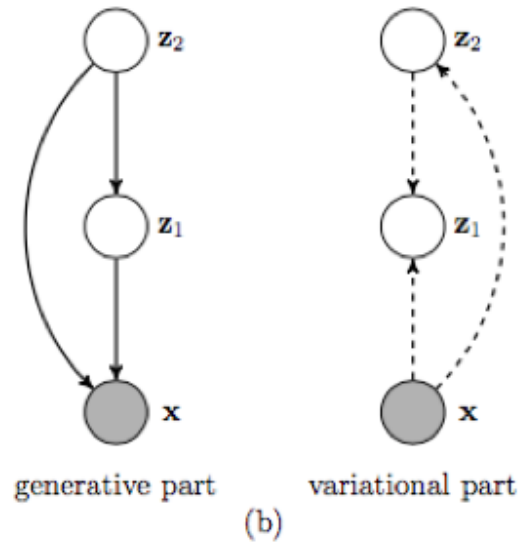
$$\therefore N \cdot \ln p(x | z) + \frac{N}{2} \cdot \sum_{j=1}^J \{ 1 + \ln \sigma_{z_{i,j}}^{2'} - \mu_{z_{i,j}}^{2'} - \sigma_{z_{i,j}}^{2'} \} + \frac{1}{2} \sum_{M=1}^M \{ 1 + \ln \sigma_{\theta_m}^{2(l)} - \mu_{\theta_m}^{2(l)} - \sigma_{\theta_m}^{2(l)} \}$$

EXPERIMENTAL RESULT

STANDARD VAE



HIERARCHICAL VAE



$$p_{\lambda}(\mathbf{z}_1|\mathbf{z}_2) = \mathcal{N}(\mathbf{z}_1|\mu_{\lambda}(\mathbf{z}_2), \text{diag}(\sigma_{\lambda}^2(\mathbf{z}_2))),$$

$$q_{\phi}(\mathbf{z}_1|\mathbf{x}, \mathbf{z}_2) = \mathcal{N}(\mathbf{z}_1|\mu_{\phi}(\mathbf{x}, \mathbf{z}_2), \text{diag}(\sigma_{\phi}^2(\mathbf{x}, \mathbf{z}_2))),$$

$$q_{\psi}(\mathbf{z}_2|\mathbf{x}) = \mathcal{N}(\mathbf{z}_2|\mu_{\psi}(\mathbf{x}), \text{diag}(\sigma_{\psi}^2(\mathbf{x})))$$

EXPERIMENTAL RESULTS

- **Table**

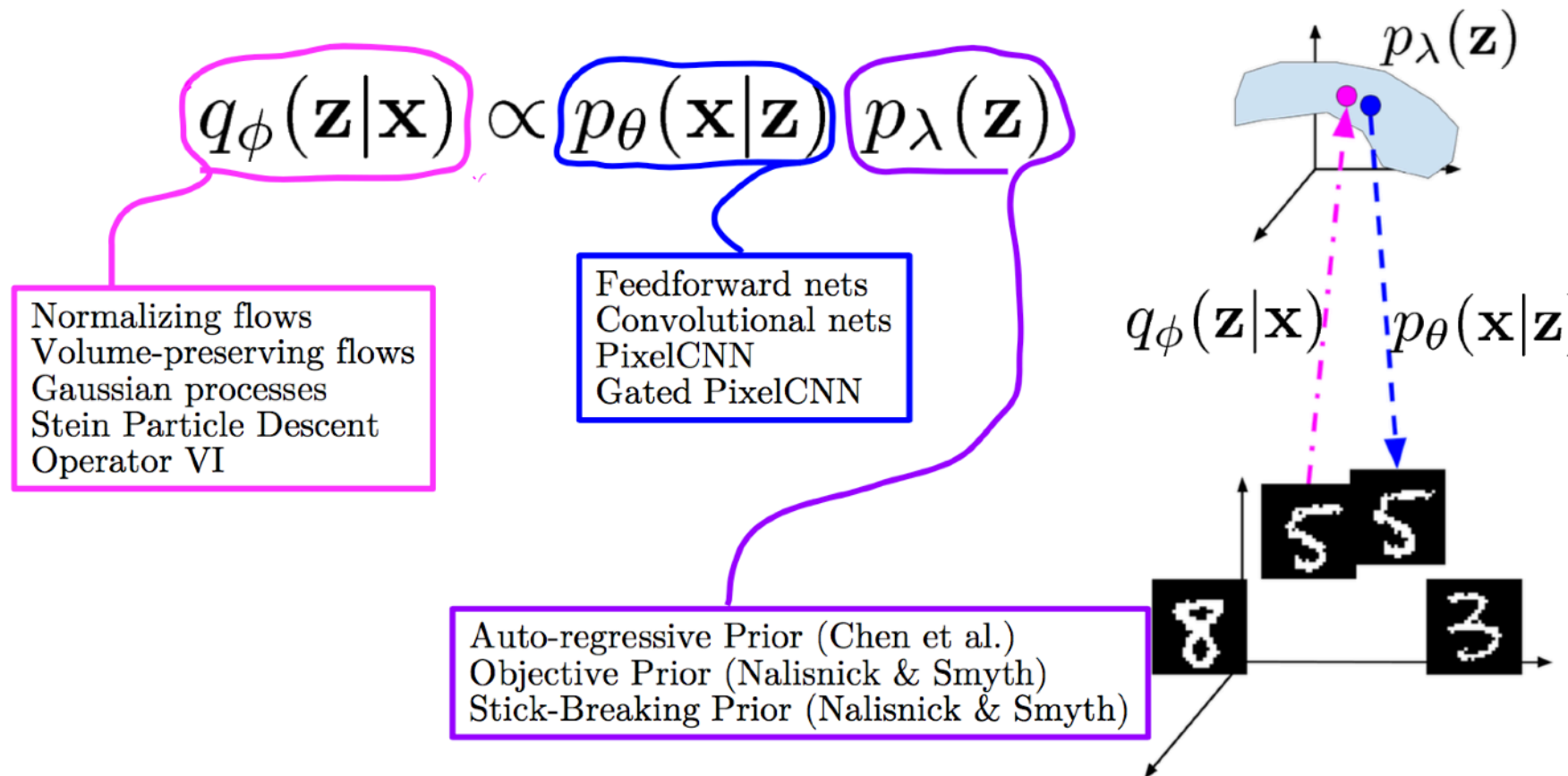
Epoch	VAE		Hierarchical VAE	
	Variational Bound	Log Marginal Likelihood	Variational Bound	Log Marginal Likelihood
1	212.87	-210.73	214.06	-214.11
50	152.83	-151.75	152.75	-152.85
100	154.24	-153.59	149.95	-150.15
150	152.32	-152.15	139.21	-139.29
200	138.41	-137.42	132.70	-132.82

WHAT'S NEXT?

CURRENT RESEARCH

RESEARCH SUMMARY

DGM: Variational Auto-Encoder



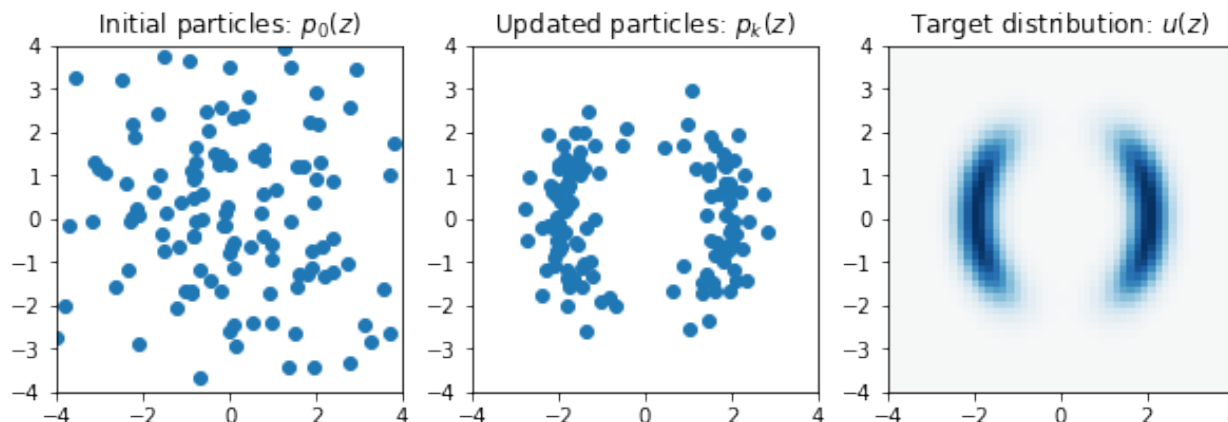
J. Tomczak (2018) VAM with a Vamp Prior at MPI Tübingen presentation

NORMALIZING FLOW / STEIN METHOD

- **Iterative update** : $z_i \sim q_\phi(z_i | x_i)$
 - Based on change of variable
 - Objective : minimizing $KLD(q || p)$
 - Normalizing flow : using tractable determinant of Jacobian
 - Stein method : using the RKHS kernel to describe direction for minimizing KLD

$$T_l^*(x) = x + \epsilon_l \phi_{q_l, p}^*(x).$$

$$q_0 \xrightarrow{T_0^*} q_1 \xrightarrow{T_1^*} q_2 \xrightarrow{T_2^*} \dots \rightarrow q_\infty = p(\phi_{q_l, p}^* = 0)$$



VAE WITH VAMPRIOR

- **Main Idea**

- . Prior Distribution = Aggregate Posterior Distribution

$$p_{\lambda}(z) = \frac{1}{K} \sum_{k=1}^K q(z|x)$$

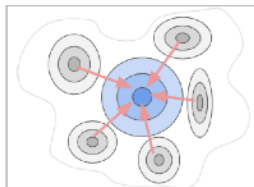
- . Pseudo inputs

- . Users should select learnable pseudo inputs.

- . If pseudo inputs are randomly chosen, we couldn't expect better performance.

- . Expectation for this approach

$$\min. \mathbb{H}[q(z)] + \textcircled{KL[q(z)||p_{\lambda}(z)]}$$



Standard prior is too strong and overregularizes the encoder.

What is the "optimal" prior?

VAE WITH VAMPRIOR

- Variational Bound

$$\begin{aligned}\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} [\ln p(\mathbf{x})] &\geq \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} [\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\ln p_\theta(\mathbf{x}|\mathbf{z})]] + \\ &\quad + \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} [\mathbb{H}[q_\phi(\mathbf{z}|\mathbf{x})]] + \\ &\quad - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [-\ln p_\lambda(\mathbf{z})]\end{aligned}$$

$$\max_{p_\lambda(\mathbf{z})} -\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [-\ln p_\lambda(\mathbf{z})] + \beta \left(\int p_\lambda(\mathbf{z}) d\mathbf{z} - 1 \right)$$



$$p_\lambda^*(\mathbf{z}) = \frac{1}{N} \sum_{n=1}^N q_\phi(\mathbf{z}|\mathbf{x}_n)$$

1 of N observations

VAE WITH VAMPRIOR

• Experimental Result

DATASET	VAE ($L = 1$)		HVAE ($L = 2$)		CONVHVAE ($L = 2$)		PIXELHVAE ($L = 2$)	
	standard	VampPrior	standard	VampPrior	standard	VampPrior	standard	VampPrior
staticMNIST	-88.56	-85.57	-86.05	-83.19	-82.41	-81.09	-80.58	-79.78
dynamicMNIST	-84.50	-82.38	-82.42	-81.24	-80.40	-79.75	-79.70	-78.45
Omniglot	-108.50	-104.75	-103.52	-101.18	-97.65	-97.56	-90.11	-89.76
Caltech 101	-123.43	-114.55	-112.08	-108.28	-106.35	-104.22	-85.51	-86.22
Frey Faces	4.63	4.57	4.61	4.51	4.49	4.45	4.43	4.38
Histopathology	6.07	6.04	5.82	5.75	5.59	5.58	4.84	4.82

Table 2: Test LL for static MNIST.

MODEL	LL
VAE ($L = 1$) + NF [32]	-85.10
VAE ($L = 2$) [6]	-87.86
IWAE ($L = 2$) [6]	-85.32
HVAE ($L = 2$) + SG	-85.89
HVAE ($L = 2$) + MoG	-85.07
HVAE ($L = 2$) + VAMPRIOR <i>data</i>	-85.71
HVAE ($L = 2$) + VAMPRIOR	-83.19
AVB + AC ($L = 1$) [28]	-80.20
VLAe [7]	-79.03
VAE + IAF [18]	-79.88
CONVHVAE ($L = 2$) + VAMPRIOR	-81.09
PIXELHVAE ($L = 2$) + VAMPRIOR	-79.78

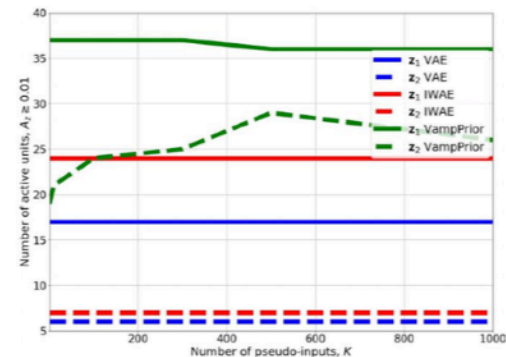


Figure 3: A comparison between two-level VAE and IWAE with the standard normal prior and their VampPrior counterpart in terms of number of active units for varying number of pseudo-inputs on static MNIST.

REFERENCES

- Kingma and Welling(2013), *Auto-encoding Variational Bayes*, ICLR2014
- Bishop(2006), Pattern recognition and Machine learning
- Rezende, Mohamed and Wiestra(2014) *Stochastic Backpropagation and Approximate Inference in Deep Generative Models*, ICML2014
- Liu and Wang(2016), *Stein Variational Gradient Descent: A General Purpose Bayesian Inference Algorithm*, NIPS2016
- Tomczak and Welling(2018), *VAE with a VampPrior*, AISTAT2018

J. Tomczak (2018) VAM with a Vamp Prior at MPI Tübingen presentation



ANY QUESTIONS?