



电子科技大学
University of Electronic Science and Technology of China



Neural Variational Inference

Di Wu



Data Mining Lab, Big Data Research Center, UESTC
Email: wudi.araragi@qq.com
Homepage: <http://dm.uestc.edu.cn>

Begin with VAE

Variational auto-encoder is used to perform approximate inference on probabilistic models which have intractable posterior distribution over latent variables and parameters. It fits an approximate inference model (also called recognition model, just an encoder) to the true posterior using an estimator of the ELBO.

Two key points:

- ◆ reparameterization trick
- ◆ efficient optimization of ELBO by stochastic gradient

VAE Model

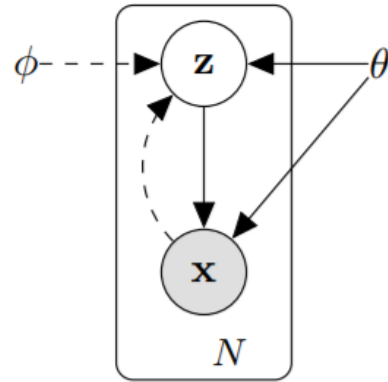


Figure 1: The type of directed graphical model under consideration. Solid lines denote the generative model $p_{\theta}(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z})$, dashed lines denote the variational approximation $q_{\phi}(\mathbf{z}|\mathbf{x})$ to the intractable posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$. The variational parameters ϕ are learned jointly with the generative model parameters θ .

Marginal log likelihood is $\log p_{\theta}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)})$

Rewrite the likelihood using a variational distribution $q_{\phi}(\mathbf{z}|\mathbf{x})$:

$$\log p_{\theta}(\mathbf{x}^{(i)}) = D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z}|\mathbf{x}^{(i)})) + \mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})$$

$$\log p_{\theta}(\mathbf{x}^{(i)}) \geq \mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [-\log q_{\phi}(\mathbf{z}|\mathbf{x}) + \log p_{\theta}(\mathbf{x}, \mathbf{z})]$$

$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = -D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z})) + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})] \quad (*1*)$$

Reparameterization Trick: Estimator of the ELBO

We want to differentiate and optimize the lower bound \mathcal{L} in (1) w.r.t ϕ and θ , the main difficulty is the gradient of ϕ . With a well chosen posterior $q_\phi(\mathbf{z}|\mathbf{x})$, we can reparameterize variable $\tilde{\mathbf{z}} \sim q_\phi(\mathbf{z}|\mathbf{x})$ using a differentiable transformation $g_\phi(\epsilon, \mathbf{x})$, ϵ is an auxiliary noise variable:

$$\tilde{\mathbf{z}} = g_\phi(\epsilon, \mathbf{x}) \quad \text{with} \quad \epsilon \sim p(\epsilon)$$

Monte Carlo estimates:

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [f(\mathbf{z})] = \mathbb{E}_{p(\epsilon)} [f(g_\phi(\epsilon, \mathbf{x}^{(i)}))] \simeq \frac{1}{L} \sum_{l=1}^L f(g_\phi(\epsilon^{(l)}, \mathbf{x}^{(i)})) \quad \text{where} \quad \epsilon^{(l)} \sim p(\epsilon)$$

The ELBO can be rewritten as:



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

$$\text{where} \quad \mathbf{z}^{(i,l)} = g_\phi(\epsilon^{(i,l)}, \mathbf{x}^{(i)}) \quad \text{and} \quad \epsilon^{(l)} \sim p(\epsilon)$$

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))

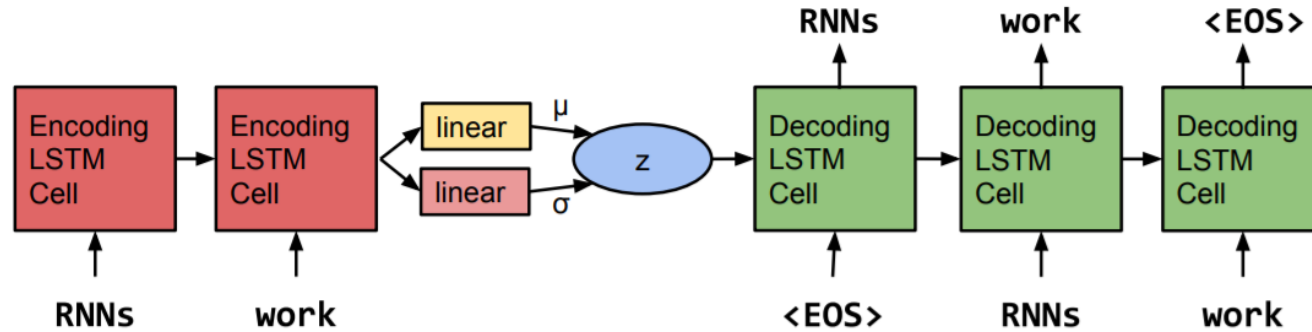
$\theta, \phi \leftarrow$ Update parameters using gradients \mathbf{g} (e.g. SGD or Adagrad [DHS10])

until convergence of parameters (θ, ϕ)

return θ, ϕ

Flexibility of the choice of the prior and the design of variational posterior

How to Apply VAE Framework to NLP



Simple approach:

$$\begin{aligned} \text{prior} \quad p(z) &= \mathcal{N}(\mu_0, \sigma_0^2) \\ \text{posterior} \quad q(z|x) &= \mathcal{N}(\mu = f(x), \sigma = g(x)) \end{aligned}$$

Figure 1: The core structure of our variational autoencoder language model. Words are represented using a learned randomly-initialized dictionary of embedding vectors. \vec{z} is a vector-valued latent variable with a Gaussian prior and an approximate posterior parameterized by the encoder's outputs μ and σ . $\langle \text{EOS} \rangle$ marks the end of each sequence.

Variable Z can be seen as the sentence semantic(global feature, like topic)

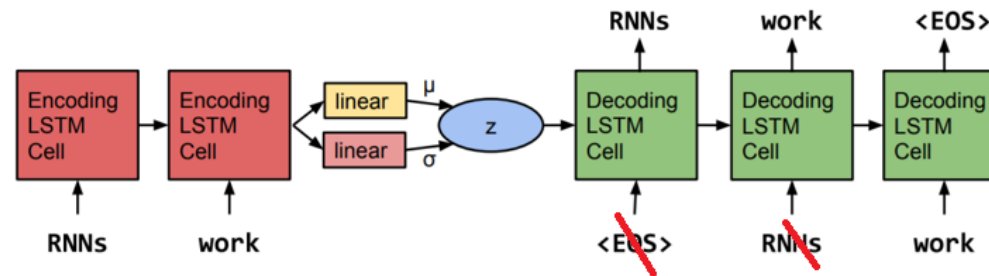
It seems okay, but:

$$\mathcal{L}(\theta; x) = -\text{KL}(q_{\theta}(\vec{z}|x) || p(\vec{z})) + \mathbb{E}_{q_{\theta}(\vec{z}|x)} [\log p_{\theta}(x|\vec{z})]$$

Because RNN can express arbitrary distributions over the output sentences, so RNN can achieve optimal likelihood even without Z, so KL will fall down to zero, actually Z doesn't be learned

How to alleviate:

- Word dropout:



- KL cost annealing: $\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})] - W^* D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)}) || p_{\theta}(\mathbf{z}))$

W increases gradually from 0

Neural Variational Inference for Text Processing(1)

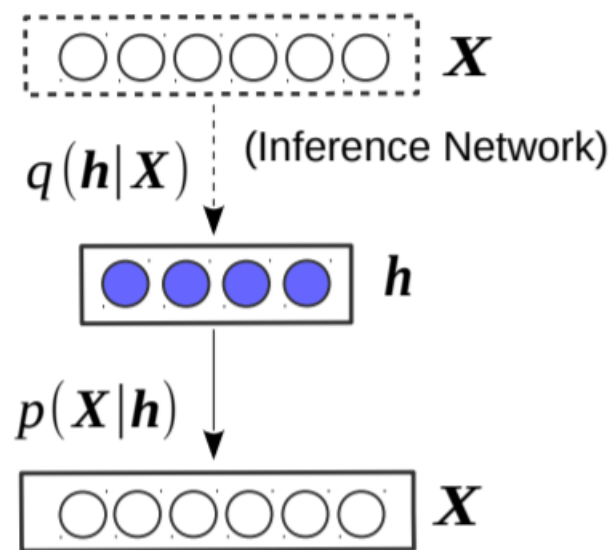


Figure 1. NVDM for document modelling.

都是套路：

\mathbf{X} : document representation (e.g. bag of words), x_i is the i th word (one-hot)

\mathbf{h} : continuous hidden variable which generates all the words independently

prior $p_\theta(\mathbf{h})$ is a Gaussian

$$p_\theta(x_i|\mathbf{h}) = \frac{\exp\{E(x_i; \mathbf{h}, \theta)\}}{\sum_{j=1}^{|\mathcal{V}|} \exp\{E(x_j; \mathbf{h}, \theta)\}}, \text{ where } E(x_i; \mathbf{h}, \theta) = \mathbf{h}^T \mathbf{R} x_i - b_{x_i}, \mathbf{R} \in \mathbb{R}^{K \times |\mathcal{V}|}$$

posterior $q_\phi(\mathbf{h}|\mathbf{X}) = \mathcal{N}(\mathbf{h}|\mu(\mathbf{X}), \text{diag}(\sigma^2(\mathbf{X})))$

$$\pi = g(f_X^{MLP}(\mathbf{X}))$$

$$\mu = l_1(\pi), \log \sigma = l_2(\pi)$$

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

Neural Variational Inference for Text Processing(2)

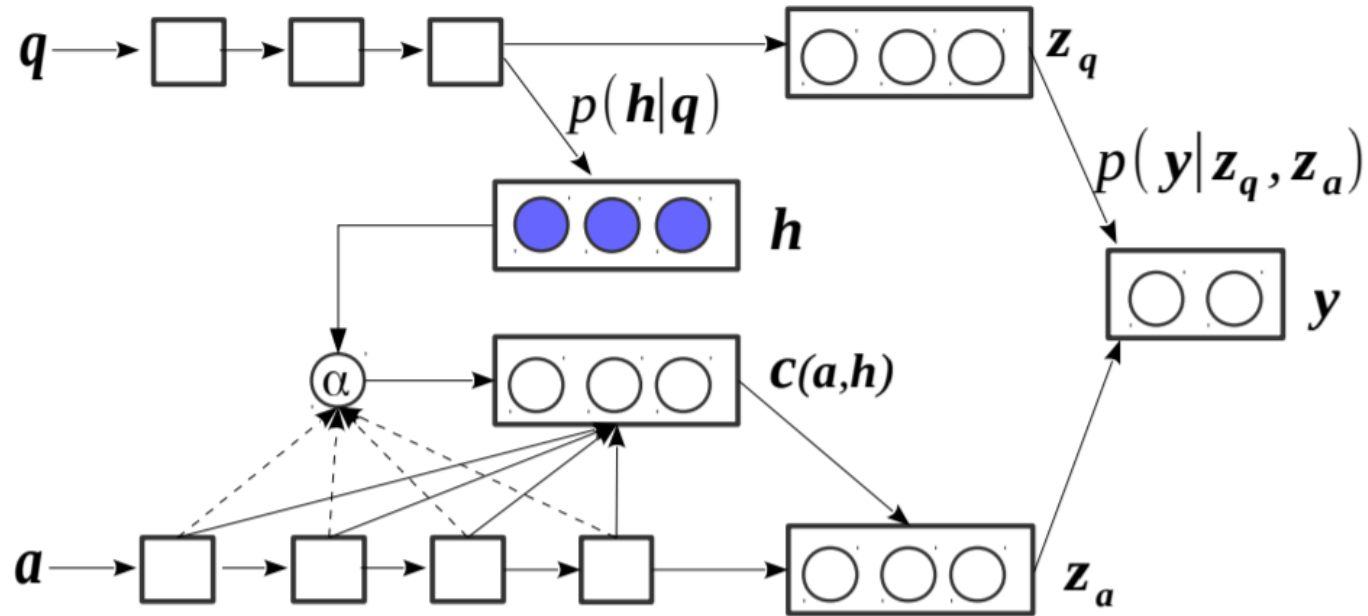


Figure 2. NASM for question answer selection.

Scenario:

Given a question q , a set of candidate answer (a_1, a_2, \dots, a_n) and judgment (y_1, y_2, \dots, y_n) where $y_m = 1$ if a_m is the answer. so each train data point is the triple (p, a, y)

Neural Variational Inference for Text Processing(2)

Model specification :

Prior :

$$p_{\theta}(\mathbf{h}|\mathbf{q}) = \mathcal{N}(\mathbf{h}|\boldsymbol{\mu}(\mathbf{q}), \text{diag}(\boldsymbol{\sigma}^2(\mathbf{q})))$$

$$\boldsymbol{\pi}_{\theta} = g_{\theta}(f_q^{\text{LSTM}}(\mathbf{q})) = g_{\theta}(\mathbf{s}_q(|\mathbf{q}|))$$

$$\boldsymbol{\mu}_{\theta} = l_1(\boldsymbol{\pi}_{\theta}), \log \boldsymbol{\sigma}_{\theta} = l_2(\boldsymbol{\pi}_{\theta})$$

Variational posterior :

$$q_{\phi}(\mathbf{h}|\mathbf{q}, \mathbf{a}, \mathbf{y}) = \mathcal{N}(\mathbf{h}|\boldsymbol{\mu}_{\phi}(\mathbf{q}, \mathbf{a}, \mathbf{y}), \text{diag}(\boldsymbol{\sigma}_{\phi}^2(\mathbf{q}, \mathbf{a}, \mathbf{y})))$$

$$\boldsymbol{\pi}_{\phi} = g_{\phi}(f_q^{\text{LSTM}}(\mathbf{q}), f_a^{\text{LSTM}}(\mathbf{a}), f_y(\mathbf{y}))$$

$$= g_{\phi}(\mathbf{s}_q(|\mathbf{q}|), \mathbf{s}_a(|\mathbf{a}|), \mathbf{s}_y)$$

$$\boldsymbol{\mu}_{\phi} = l_3(\boldsymbol{\pi}_{\phi}), \log \boldsymbol{\sigma}_{\phi} = l_4(\boldsymbol{\pi}_{\phi})$$

Generative process:

$$\alpha(i) \propto \exp(\mathbf{W}_{\alpha}^T \tanh(\mathbf{W}_h \mathbf{h} + \mathbf{W}_s \mathbf{s}_a(i)))$$

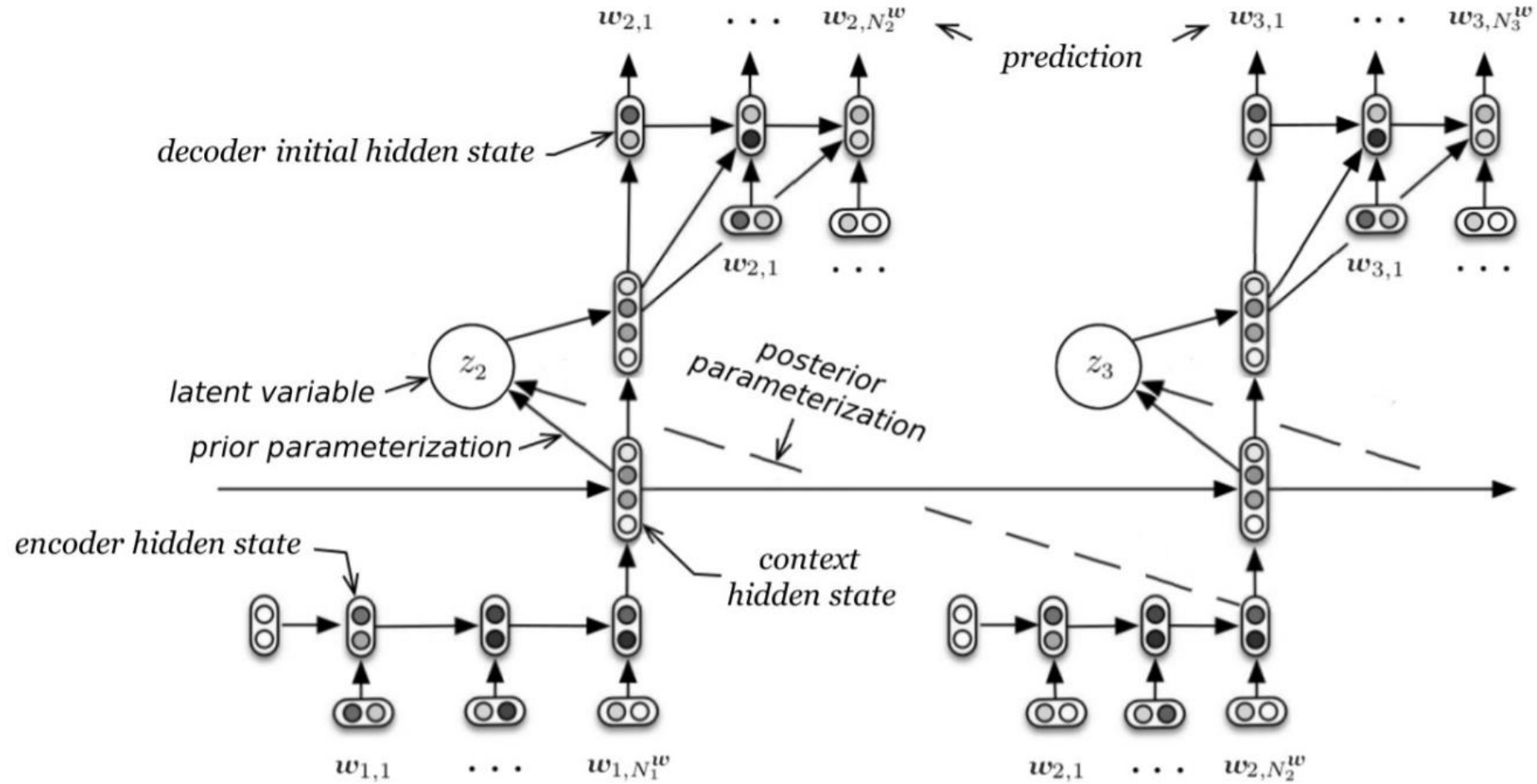
$$\mathbf{c}(\mathbf{a}, \mathbf{h}) = \sum_i \mathbf{s}_a(i) \alpha(i)$$

$$\mathbf{z}_a(\mathbf{a}, \mathbf{h}) = \tanh(\mathbf{W}_a \mathbf{c}(\mathbf{a}, \mathbf{h}) + \mathbf{W}_n \mathbf{s}_a(|\mathbf{a}|))$$

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

$$\text{ELBO} : \quad \mathcal{L} = \mathbb{E}_{q_{\phi}(\mathbf{h})}[\log p_{\theta}(\mathbf{y} | \mathbf{z}_q(\mathbf{q}), \mathbf{z}_a(\mathbf{a}, \mathbf{h}))] - D_{\text{KL}}(q_{\phi}(\mathbf{h}) || p_{\theta}(\mathbf{h} | \mathbf{q}))$$

Neural Variational Inference for generating dialogues



Neural Variational Inference for generating dialogues

Prior :
$$P_{\theta}(\mathbf{z}_n \mid \mathbf{w}_1, \dots, \mathbf{w}_{n-1}) = \mathcal{N}(\boldsymbol{\mu}_{\text{prior}}(\mathbf{w}_1, \dots, \mathbf{w}_{n-1}), \Sigma_{\text{prior}}(\mathbf{w}_1, \dots, \mathbf{w}_{n-1})),$$

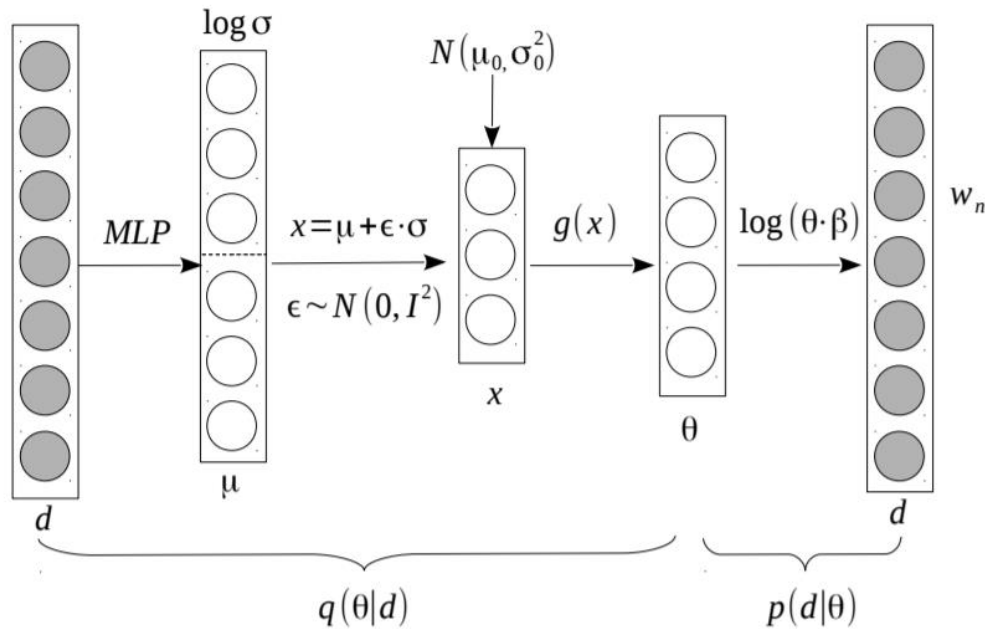
Generative process:
$$P_{\theta}(\mathbf{w}_n \mid \mathbf{z}_n, \mathbf{w}_1, \dots, \mathbf{w}_{n-1}) = \prod_{m=1}^{M_n} P_{\theta}(w_{n,m} \mid \mathbf{z}_n, \mathbf{w}_1, \dots, \mathbf{w}_{n-1}, w_{n,1}, \dots, w_{n,m-1})$$

Variational posterior :
$$Q_{\psi}(\mathbf{z}_n \mid \mathbf{w}_1, \dots, \mathbf{w}_n) = \mathcal{N}(\boldsymbol{\mu}_{\text{posterior}}(\mathbf{w}_1, \dots, \mathbf{w}_n), \Sigma_{\text{posterior}}(\mathbf{w}_1, \dots, \mathbf{w}_n))$$

ELBO :
$$\log P_{\theta}(\mathbf{w}_1, \dots, \mathbf{w}_N) \geq \sum_{n=1}^N -\text{KL} [Q_{\psi}(\mathbf{z}_n \mid \mathbf{w}_1, \dots, \mathbf{w}_n) \parallel P_{\theta}(\mathbf{z}_n \mid \mathbf{w}_1, \dots, \mathbf{w}_{n-1})] \\ + \mathbb{E}_{Q_{\psi}(\mathbf{z}_n \mid \mathbf{w}_1, \dots, \mathbf{w}_n)} [\log P_{\theta}(\mathbf{w}_n \mid \mathbf{z}_n, \mathbf{w}_1, \dots, \mathbf{w}_{n-1})],$$

Neural Variational Topic Model

Document topic distribution is a multinomial(discrete),so just transform the Gaussian variable by softmax, the rest is the same.



$g(x)$ is the transform function:

$$x \sim \mathcal{N}(\mu_0, \sigma_0^2)$$

$$\theta = \text{softmax}(W_1^T x)$$

ELBO:

$$\mathcal{L}_d = \mathbb{E}_{q(\theta|d)} \left[\sum_{n=1}^N \log \sum_{z_n} [p(w_n | \beta_{z_n}) p(z_n | \theta)] \right]$$


$$- D_{KL} [q(\theta|d) || p(\theta | \mu_0, \sigma_0^2)]$$

Figure 3. Network structure of the inference model $q(\theta | d)$, and of the generative model $p(d | \theta)$.

Neural Variational Topic Model(non-parametric version)

Stick Breaking Process

$$\nu_k \sim \text{Beta}(1, \alpha) \quad \pi_k = \nu_k \prod_{l=1}^{k-1} (1 - \nu_l) = \nu_k \left(1 - \sum_{l=1}^{k-1} \pi_l\right)$$



$$\boldsymbol{\pi} = \{\pi_k\}_{k=1}^{\infty} \quad 0 \leq \pi_k \leq 1 \text{ and } \sum_{k=1}^{\infty} \pi_k = 1$$

Kumaraswamy distribution(similar to Beta distribution,more suitable for reparameterization trick)

$$\text{Kumaraswamy}(x; a, b) = abx^{a-1}(1 - x^a)^{b-1}$$

Inverse CDF: $x = (1 - u^{\frac{1}{b}})^{\frac{1}{a}}$, where $u \sim \text{Uniform}(0, 1)$

Neural Variational Topic Model(non-parametric version)

Generative Story:

Stick breaking process

- Draw a topic distribution $\pi \sim \text{GEM}(\alpha)$
- Then we get a distribution $G(\theta; \pi, \Theta) = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}(\theta)$
- For each word w_i in the document: 1) draw a topic $\hat{\theta}_i \sim G(\theta; \pi, \Theta)$; 2) $w_i \sim \text{Cat}(\hat{\theta}_i)$

Here, we want to approximate the posterior distribution of $\nu_k \sim \text{Beta}(1, \alpha)$

Prior: $p(\boldsymbol{\nu}|\alpha)$ is products of $K - 1$ Beta(1, α)

Variational posterior: $[a_1, \dots, a_{K-1}; b_1, \dots, b_{K-1}] = g(\mathbf{w}_{1:N}; \psi)$

$$q_\psi(\boldsymbol{\nu}|\mathbf{w}_{1:N}) = \prod_{k=1}^{K-1} \kappa(\nu_k; a_k, b_k)$$

Likelihood: $\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_{K-1}, \pi_K) = \left(\nu_1, \nu_2(1 - \nu_1), \dots, \nu_{k-1} \prod_{l=1}^{K-2} (1 - \nu_l), \prod_{l=1}^{K-1} (1 - \nu_l) \right)$

$$p(\mathbf{w}_{1:N}, \boldsymbol{\pi}, \hat{\boldsymbol{\theta}}_{1:N}|\alpha, \Theta) = p(\boldsymbol{\pi}|\alpha) \prod_{i=1}^N p(w_i|\hat{\boldsymbol{\theta}}_i)p(\hat{\boldsymbol{\theta}}_i|\boldsymbol{\pi}, \Theta)$$

$$p(\mathbf{w}_{1:N}, \boldsymbol{\pi}|\alpha, \Theta) = p(\boldsymbol{\pi}|\alpha) \prod_{i=1}^N p(w_i|\boldsymbol{\pi}, \Theta)$$

ELBO: $\mathcal{L}(\mathbf{w}_{1:N}|\Phi, \psi) = \mathbb{E}_{q_\psi(\boldsymbol{\nu}|\mathbf{w}_{1:N})} [\log p(\mathbf{w}_{1:N}|\boldsymbol{\pi}, \Phi)] - \text{KL}(q_\psi(\boldsymbol{\nu}|\mathbf{w}_{1:N})||p(\boldsymbol{\nu}|\alpha))$

References:

1. **Auto-Encoding Variational Bayes**, ICLR 2014
2. **Generating Sentences From a Continuous Spaces**, ICLR 2016
3. **Neural Variational Inference for Text Processing**, ICML 2016
4. **A Hierarchical Latent Variable Encoder-Decoder Model for Generating Dialogues**, AAAI 2017
5. **Discovering Discrete Latent Topics with Neural Variational Inference**, ICML 2017
6. **A Bayesian Nonparametric Topic Model with Variational Auto-Encoders**, ICLR 2018 under review
7. **A Conditional Variational Framework for Dialog Generation**, ACL 2017