

## [ Paper review 41 ]

---

# Neural Variational Inference for Text Processing

---

( Miao, et al., 2016 )

## [ Contents ]

---

1. Abstract
2. Introduction
3. Neural Variational Inference Framework
4. NVDM (Neural Variational Document Model)

## 1. Abstract

---

Introduce **Variational Inference** framework for **generative & conditional models of text**

- traditional) analytic approximation for intractable distn
- this paper) construct **inference network** conditioned on the discrete text input

Validate using 2 applications

- 1) **generative document modeling**
- 2) **supervised question answering**

## 2. Introduction

---

Probabilistic generative models in NLP

- 1) ability use unlabelled data effectively
- 2) incorporate abundant linguistic features
- 3) learn interpretable dependencies among data

→ but as the model becomes complex.... becomes intractable! ( due to high dimensional integrals )

→ use MCMC or VI

### Problem of MCMC & VI

- MCMC : computational cost
- VI : confined due to underestimation of posterior variance

Introduces **NEURAL VARIATIONAL FRAMEWORK** for generative models of text, inspired by **VAE**

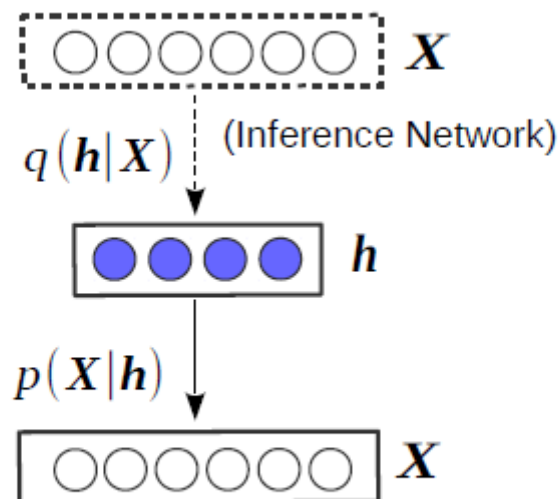
→ main idea : building **inference network** with DNN

- analytic approximation ( $X$ )  
model the posterior probability ( $O$ ) .... thus strong generalization ability
- due to DNN : capture of learning complex non-linear distn
- use reparam tricks  
( to train by back-prop of unbiased & low variance gradients w.r.t latent variables )

Propose

- NVDM (Neural Variational Document Model)
- NASM (Neural Answer Selection Model)

## NVDM



*Figure 1. NVDM for document modelling.*

- **unsupervised** generative model of text
- extract "continuous semantic latent variable" for each document
- like VAE
  - **MLP encoder** (  $X$  : bag-of-words documents  $\rightarrow Z$  : continuous latent distn )
  - **Softmax decoder** ( reconstructs document! )  
( each word is generated directly from dense continuous document representation )

## NASM

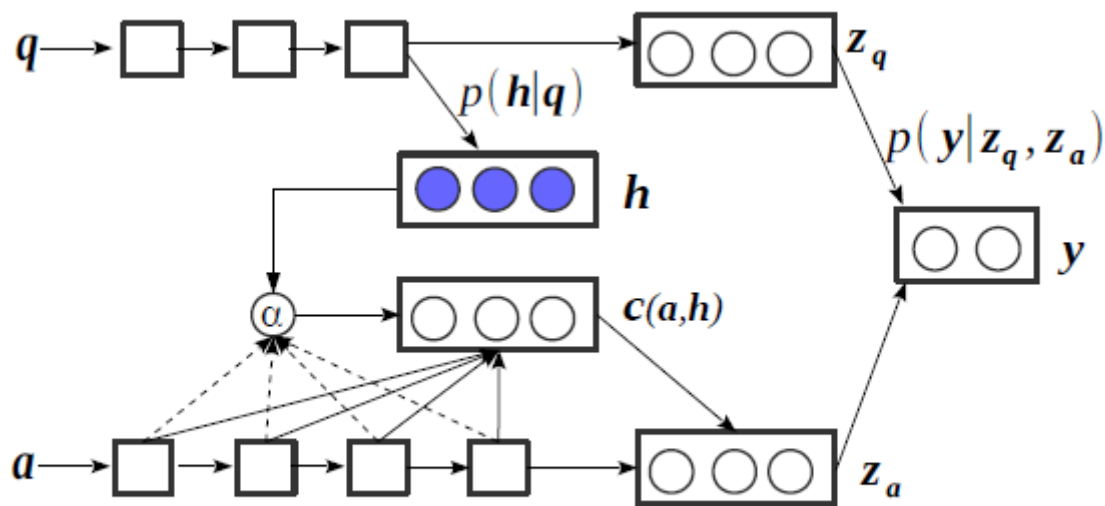


Figure 2. NASM for question answer selection.

- **supervised** conditional model which imbues **LSTMs**
- **attention model**  
( to focus on the phrases of an answer that are strongly connected to the question semantics )

Summary

demonstrate the effectiveness of **Neural Variational Inference for text processing** on 2 tasks

the models are simple, expressive, trained efficiently ( with scalable stochastic gradient back-prop )

suitable for unsupervised & supervised task

### 3. Neural Variational Inference Framework

latent variable model : popular in NLP, but hard to be effective & efficient in complex structures

→ thus propose generic "**neural variational inference framework**" that can be applied to

- unsupervised NVDM
- supervised NASM

Generative model

- latent variable  $h$   
( = stochastic units in DNN )
- $x, y$  : observed parent & child nodes of  $h$
- joint pdf :  $p_{\theta}(x, y) = \sum_h p_{\theta}(y | h) p_{\theta}(h | x) p(x)$

ELBO :

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{q(\mathbf{h})} [\log p_{\theta}(\mathbf{y} | \mathbf{h}) p_{\theta}(\mathbf{h} | \mathbf{x}) p(\mathbf{x}) - \log q(\mathbf{h})] \\ &\leq \log \int \frac{q(\mathbf{h})}{q(\mathbf{h})} p_{\theta}(\mathbf{y} | \mathbf{h}) p_{\theta}(\mathbf{h} | \mathbf{x}) p(\mathbf{x}) d\mathbf{h} = \log p_{\theta}(\mathbf{x}, \mathbf{y})\end{aligned}$$

- $q(\mathbf{h})$  should approach true posterior  $p(\mathbf{h} | \mathbf{x}, \mathbf{y})$
- parameterized diagonal Gaussian  
 $q_{\phi}(\mathbf{h} | \mathbf{x}, \mathbf{y}) = \mathcal{N}(\mathbf{h} | \boldsymbol{\mu}(\mathbf{x}, \mathbf{y}), \text{diag}(\boldsymbol{\sigma}^2(\mathbf{x}, \mathbf{y})))$ .

3 steps to construct **Inference network**

- step 1) Construct **vector representations of the observed variables** :  
 $u = f_x(x), v = f_y(y)$ .
  - $f_x(\cdot)$  and  $f_y(\cdot)$  : DNN
- step 2) Assemble a **joint representation**:  $\pi = g(u, v)$ .
  - $g(\cdot)$  : MLP that concatenates vector representations of the conditioning variables
- step 3) **Parameterise the variational distribution** over the latent variable:  
 $\mu = l_1(\pi), \log \sigma = l_2(\pi)$ 
  - $l(\cdot)$  : linear transformation & outputs the params of Gaussian

During training,

- model params  $\theta$
- inference network params  $\phi$

are updated using stochastic-backprop

## (1) model params $\theta$

$$\nabla_{\theta} \mathcal{L} \simeq \frac{1}{L} \sum_{l=1}^L \nabla_{\theta} \log p_{\theta}(\mathbf{y} | \mathbf{h}^{(l)}) p_{\theta}(\mathbf{h}^{(l)} | \mathbf{x})$$

where  $\mathbf{h} \sim q_{\phi}(\mathbf{h} | \mathbf{x}, \mathbf{y})$

## (2) inference network params $\phi$

need reparam trick!

- $\mathbf{h} = \boldsymbol{\mu} + \boldsymbol{\sigma} \cdot \boldsymbol{\epsilon}$ .  
where  $\boldsymbol{\epsilon}^{(l)} \sim \mathcal{N}(0, I)$
- to reduce variance!

update of  $\phi$  can be carried out, using gradients w.r.t.  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}$

$$s(\mathbf{h}) = \log p_{\theta}(\mathbf{y} | \mathbf{h}) p_{\theta}(\mathbf{h} | \mathbf{x}) - \log q_{\phi}(\mathbf{h} | \mathbf{x}, \mathbf{y})$$

$$\nabla_{\mu} \mathcal{L} \simeq \frac{1}{L} \sum_{l=1}^L \nabla_{\mathbf{h}^{(l)}} [s(\mathbf{h}^{(l)})].$$

$$\nabla_{\sigma} \mathcal{L} \simeq \frac{1}{2L} \sum_{l=1}^L \boldsymbol{\epsilon}^{(l)} \nabla_{\mathbf{h}^{(l)}} [s(\mathbf{h}^{(l)})].$$

## 4. NVDM (Neural Variational Document Model)

---

### "Unsupervised learning"

$h \in R^K$  : continuous hidden variable , used to represent its semantic content

interpret NVDM as VAE

- MLP encoder :  $q(\mathbf{h} \mid \mathbf{X})$
- softmax decoder :  $p(\mathbf{X} \mid \mathbf{h}) = \prod_{i=1}^N p(\mathbf{x}_i \mid \mathbf{h})$

To maximize log likelihood ( =  $\log \sum_h p(\mathbf{X} \mid \mathbf{h})p(\mathbf{h})$  ),

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

- $N$  : number of words in the document
- $p(\mathbf{h})$  : Gaussian prior for  $h$

### Encoder

- inference network
- modeled by Gaussian
- $q_\phi(\mathbf{h} \mid \mathbf{X}) = \mathcal{N}(\mathbf{h} \mid \boldsymbol{\mu}(\mathbf{X}), \text{diag}(\boldsymbol{\sigma}^2(\mathbf{X})))$ .  
where  $\boldsymbol{\pi} = g(f_X^{\text{MLP}}(\mathbf{X}))$ . and  $\boldsymbol{\mu} = l_1(\boldsymbol{\pi})$ ,  $\log \boldsymbol{\sigma} = l_2(\boldsymbol{\pi})$

### Decoder

- modeled by softmax ( = multinomial logistic regression )
- $p_\theta(\mathbf{x}_i \mid \mathbf{h}) = \frac{\exp\{-E(\mathbf{x}_i; \mathbf{h}, \theta)\}}{\sum_{j=1}^{|\mathbf{V}|} \exp\{-E(\mathbf{x}_j; \mathbf{h}, \theta)\}}$ .  
where  $E(\mathbf{x}_i; \mathbf{h}, \theta) = -\mathbf{h}^T \mathbf{R} \mathbf{x}_i - b_{x_i}$

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

by back-prop of stochastic gradients w.r.t  $\theta$  and  $\phi$

## 5. NASM (Neural Answer Selection Model)

---

after reviewing Attention mechanism...

