# CS 6501 Natural Language Processing

## Variational Auto-encoder

Yangfeng Ji

November 7, 2018

Department of Computer Science
University of Virginia

UNIVERSITY *of* VIRGINIA | ENGINEERING

# Overview

1. Variational Inference, Review

2. Variational Auto-encoder

3. Application: Neural variational document model

# Variational Inference, Review

# Latent Variable Models

$$p(x, z; \theta) = \underbrace{p(z; \theta)}_{\text{prior}} \underbrace{p(x \mid z; \theta)}_{\text{likelihood}} \tag{1}$$

Two problems:

▶ What is $\boldsymbol{\theta}$ given a collection of data point $\{\boldsymbol{x}_n\}$?

▶ What is $p(\boldsymbol{z} \mid \boldsymbol{x}; \boldsymbol{\theta})$?

One connection

▶ Computation of $p(\boldsymbol{x}; \boldsymbol{\theta}) = \sum_z p(\boldsymbol{x}, \boldsymbol{z}; \boldsymbol{\theta})$

# Ideal Cases

... where we can compute $p(\boldsymbol{x}; \boldsymbol{\theta})$, then

- ▶ We can get an estimate of $\boldsymbol{\theta}$

$$\hat{\boldsymbol{\theta}} = \arg\max_{\boldsymbol{\theta}} p(\boldsymbol{x}; \boldsymbol{\theta}) \tag{2}$$

... where we can compute $p(x; \theta)$, then

▶ We can get an estimate of $\theta$

$$\hat{\theta} = \arg\max_{\theta} p(x; \theta) \tag{2}$$

▶ Then, for a given data point $x'$, we can compute its latent variable (representation) $z$ as

$$p(z' \mid x'; \hat{\theta}) = \frac{p(x', z'; \hat{\theta})}{p(x'; \hat{\theta})} \tag{3}$$

... where we can compute $p(x; \theta)$, then

- We can get an estimate of $\theta$

$$\hat{\theta} = \arg \max_{\theta} p(x; \theta) \qquad (2)$$

- Then, for a given data point $x'$, we can compute its latent variable (representation) $z$ as

$$p(z' \mid x'; \hat{\theta}) = \frac{p(x', z'; \hat{\theta})}{p(x'; \hat{\theta})} \qquad (3)$$

$p(x; \theta) = \sum_z p(x, z; \theta)$ is challenging in practice

# Variational Inference

▶ Approximate $p(z \mid x; \boldsymbol{\theta})$ with a variational distribution $q(z; \boldsymbol{\phi})$

$$q(z; \boldsymbol{\theta}) \approx p(z \mid x; \boldsymbol{\theta})) \tag{4}$$

# Variational Inference

▶ Approximate $p(z \mid x; \boldsymbol{\theta})$ with a variational distribution $q(z; \boldsymbol{\phi})$

$$q(z; \boldsymbol{\theta}) \approx p(z \mid x; \boldsymbol{\theta}))\qquad(4)$$

▶ Approximation by minimizing the KL divergence (still needs $p(x; \boldsymbol{\theta})$, not directly solvable)

$$\mathrm{KL}(q\|p) = E_{q(z;\boldsymbol{\phi})}[q(z; \boldsymbol{\phi})] - E_{q(z;\boldsymbol{\phi})}[p(z \mid x; \boldsymbol{\theta})]\qquad(5)$$

# Variational Inference

- Approximate $p(z \mid x; \boldsymbol{\theta})$ with a variational distribution $q(z; \boldsymbol{\phi})$

$$q(z; \boldsymbol{\theta}) \approx p(z \mid x; \boldsymbol{\theta}))\qquad(4)$$

- Approximation by minimizing the KL divergence (still needs $p(x; \boldsymbol{\theta})$, not directly solvable)

$$\text{KL}(q\|p) = E_{q(z;\boldsymbol{\phi})}[q(z; \boldsymbol{\phi})] - E_{q(z;\boldsymbol{\phi})}[p(z \mid x; \boldsymbol{\theta})]\qquad(5)$$

- or, by maximizing the evidence lower bound (solvable with some constriants)

$$\text{ELBo}(\boldsymbol{\theta}, \boldsymbol{\phi}) = E_q[p(x, z; \boldsymbol{\theta})] - E_q[q(z; \boldsymbol{\phi})]\qquad(6)$$

# Example: Mixture of Gaussians

▶ Prior:

$$p(z) = \prod_{k=1}^{K} \pi_k^{z_k}. \tag{7}$$

▶ Likelihood:

$$p(x \mid z) = \sum_{k=1}^{K} \mathcal{N}(x \mid \mu_k, \Sigma_k)^{z_k} \tag{8}$$

# Example: Mixture of Gaussians

▶ Prior:

$$p(\mathbf{z}) = \prod_{k=1}^{K} \pi_k^{z_k}. \tag{7}$$

▶ Likelihood:

$$p(\mathbf{x} \mid \mathbf{z}) = \sum_{k=1}^{K} \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \tag{8}$$
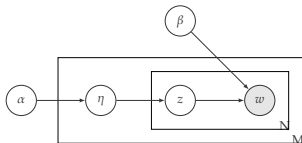
▶ $p(\mathbf{x})$ is tractable

$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \tag{9}$$

So, we can use the EM algorithm.

# Example: Latent Dirichlet Allocation

1. Choose $\eta \sim \text{Dirichlet}(\alpha)$
2. For each word $w_n$
   2.1 Choose a topic $z_n \sim \text{Multinomial}(\eta)$
   2.2 Choose a word $w_n \sim p(w_n \mid z_n; \beta)$, a multinomial probability conditioned on the topic $z_n$.



Variational inference: mean-field approximation

Choose your $p(x \mid z; \theta)$ carefully

- ▶ Categorical (in GMMs)
- ▶ Multinomial (in LDA)
- ▶ Dirichlet (in LDA)

## Question

What if $p(x \mid z; \theta)$ is not a common probability distribution?

# Variational Auto-encoder
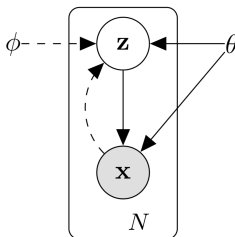
# Extension on Modeling

$$p(x, z; \theta) = p(z; \theta) \cdot p(x \mid z; \theta). \qquad (10)$$

▶ $p(x \mid z; \theta)$ is modeled by a neural network with some randomness, therefore $E_q[p(x \mid z; \theta)]$ is intractable

▶ variational distribution $q(z \mid x; \phi)$

# Graphical Representation



- $q(z \mid x; \phi)$: probabilistic encoder
- $p(x \mid z; \theta)$: probabilistic decoder

# Evidence Lower Bound

$$
\begin{aligned}
\text{ELBo}(\boldsymbol{\theta}, \boldsymbol{\phi}) \quad = \quad & E_{q(z;\phi)}[\log p(\boldsymbol{x}, \boldsymbol{z}; \boldsymbol{\theta})] - E_{q(z;\phi)}[\log q(\boldsymbol{z}; \boldsymbol{\phi})] \\
= \quad & E_{q(z;\phi)}[\log p(\boldsymbol{x} \mid \boldsymbol{z}; \boldsymbol{\theta})] + E_{q(z;\phi)}[p(\boldsymbol{z}; \boldsymbol{\theta})] \qquad (11) \\
& -E_{q(z;\phi)}[\log q(\boldsymbol{z}; \boldsymbol{\phi})] \\
= \quad & E_{q(z;\phi)}[\log p(\boldsymbol{x} \mid \boldsymbol{z}; \boldsymbol{\theta})] - \text{KL}(q(\boldsymbol{z} \mid \boldsymbol{x}; \boldsymbol{\phi}) \| p(\boldsymbol{z}; \boldsymbol{\theta})) \qquad (12)
\end{aligned}
$$

$$\nabla_{\phi} E_{q(z;\phi)}[\log p(x \mid z; \theta)] = \nabla_{\phi} \sum_z q(z; \phi) \log p(x \mid z; \theta)$$

---

1

# Gradient based Learning

$$\nabla_\phi E_{q(z;\phi)}[\log p(x \mid z; \theta)] = \nabla_\phi \sum_z q(z; \phi) \log p(x \mid z; \theta)$$

$$= \sum_z \nabla_\phi q(z; \phi) \log p(x \mid z; \theta)$$

1

# Gradient based Learning

$$\nabla_{\phi} E_{q(z;\phi)}[\log p(x \mid z; \theta)] \;=\; \nabla_{\phi} \sum_{z} q(z;\phi) \log p(x \mid z; \theta)$$

$$=\; \sum_{z} \nabla_{\phi} q(z;\phi) \log p(x \mid z; \theta)$$

$$=\; \sum_{z} q(z;\phi) \frac{\nabla_{\phi} q(z;\phi)}{q(z;\phi)} \log p(x \mid z; \theta)$$

---

1

# Gradient based Learning

$$\nabla_{\phi} E_{q(z;\phi)}[\log p(x \mid z; \theta)] = \nabla_{\phi} \sum_z q(z; \phi) \log p(x \mid z; \theta)$$

$$= \sum_z \nabla_{\phi} q(z; \phi) \log p(x \mid z; \theta)$$

$$= \sum_z q(z; \phi) \frac{\nabla_{\phi} q(z; \phi)}{q(z; \phi)} \log p(x \mid z; \theta)$$

$$= \sum_z q(z; \phi) \nabla_{\phi} \log q(z; \phi) \log p(x \mid z; \theta)[1]$$

---

[1]likelihood ratio trick, also used in policy gradient

# Gradient based Learning

$$
\begin{aligned}
\nabla_{\phi} E_{q(z;\phi)}[\log p(x \mid z; \theta)] &= \nabla_{\phi} \sum_{z} q(z; \phi) \log p(x \mid z; \theta) \\
&= \sum_{z} \nabla_{\phi} q(z; \phi) \log p(x \mid z; \theta) \\
&= \sum_{z} q(z; \phi) \frac{\nabla_{\phi} q(z; \phi)}{q(z; \phi)} \log p(x \mid z; \theta) \\
&= \sum_{z} q(z; \phi) \nabla_{\phi} \log q(z; \phi) \log p(x \mid z; \theta)^{1} \\
&\approx \frac{1}{L} \sum_{l=1}^{L} \nabla_{\phi} \log q(z^{(l)}; \phi) \log p(x \mid z^{(l)}; \theta)
\end{aligned}
$$

where $z^{(l)} \sim q(z; \phi)$

[1]likelihood ratio trick, also used in policy gradient

14

# Variance of MC Estimator

$$\nabla_\phi E_{q(z;\phi)}[\log p(x \mid z; \theta)] \approx \frac{1}{L} \sum_{l=1}^{L} \nabla_\phi \log q(z^{(l)}; \phi) \log p(x \mid z^{(l)}; \theta)$$

$$(13)$$

where $z^{(l)} \sim q(z; \phi)$.

Problems

▶ Sampling from $q(z; \phi)$ is not easy
▶ Large variance

# Reparameterization Trick

Representing $q(z; \phi)$ as a composition of a function $g$ and a simple random variable $\epsilon$

$$z = g_\phi(\epsilon, x) \qquad (14)$$

Advantages

- ▶ $g$ can be an arbitrary function (e.g., neural networks)
- ▶ $\epsilon$ is easier to sample

# Reparameterization Trick: Example

- Original form $z \sim \mathcal{N}(\mu, \sigma^2)$
- Reparameterization form:
  - $z = \mu + \sigma\epsilon$
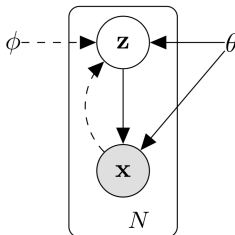  - $\epsilon \sim \mathcal{N}(0, 1)$

# Approximation of ELBo

Based on the reparameterization trick, the ELBo can be approxiamted as

$$\tilde{\mathcal{L}} = \frac{1}{L} \sum_{l=1}^{L} \log p(\boldsymbol{x}_i \mid \boldsymbol{z}_i^{(l)}) - \mathrm{KL}(q(\boldsymbol{z} \mid \boldsymbol{x}; \boldsymbol{\phi}) \| p(\boldsymbol{z}; \boldsymbol{\theta})) \quad (15)$$

where $\boldsymbol{z}_i^{(l)} = g_{\boldsymbol{\phi}}(\boldsymbol{\epsilon}_i^{(l)}, \boldsymbol{x}_i)$ and $\boldsymbol{\epsilon}_i^{(l)} \sim p(\boldsymbol{\epsilon})$.

# Comments



Encoder

- $z = g_\phi(\epsilon, x)$ with $\epsilon \sim \mathcal{N}(0, I)$
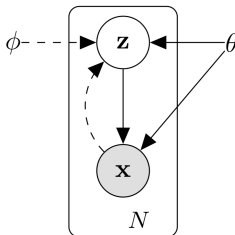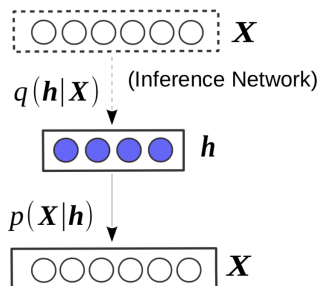- Neural networks with some simple randomness

# Comments



Encoder

- ▶ $z = g_\phi(\epsilon, x)$ with $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$
- ▶ Neural networks with some simple randomness

Decoder

- ▶ A simple nonlinear function

# Application: Neural variational document model

# Neural Variational Document Model



- an MLP document encoder $q(h \mid x; \theta)$
- a softmax decoder $p(x \mid h; \phi) = \prod_i p(x_i \mid h; \phi)$

[Miao et al., 2016]

The decoder $p(x_i \mid \boldsymbol{h})$ is defined as

$$p(x_i \mid \boldsymbol{h}; \boldsymbol{\theta}) \propto \exp(\boldsymbol{h}^\top \mathbf{R} \boldsymbol{v}_{x_i} + \boldsymbol{b}_{x_i}) \qquad (16)$$

where $\mathbf{R}$ and $\{\boldsymbol{b}_{x_i}\}$ are part of $\boldsymbol{\theta}$, and $\boldsymbol{v}_{x_i}$ is the embedding of word $x_i$.

[Miao et al., 2016]

The encoder $q(h \mid x; \phi)$ is defined as

$$q(h \mid x; \phi) = \mathcal{N}(h \mid \mu(x), \text{diag}(\sigma^2(x))) \qquad (17)$$

with

$$\mu = l_1(\pi) \qquad (18)$$
$$\log \sigma = l_2(\pi) \qquad (19)$$
$$\pi = g(x) \qquad (20)$$

where $g(\cdot)$ is a neural network and both $l_1$ and $l_2$ are linear transformations.

[Miao et al., 2016]

With the reparameterization trick, we have $h \sim q(h \mid x; \phi)$

$$h = \mu + \sigma \circ \epsilon \qquad (21)$$

where $\epsilon \sim \mathcal{N}(0, I)$.

# ELBo

The problem can be solved by directly applying the algorithm proposed in [Kingma and Welling, 2014] to maximize the following ELBo (also, Eq. 5 in [Miao et al., 2016])

$$\mathcal{L} = E[\sum_i p(x_i \mid h; \boldsymbol{\theta})] - \mathrm{KL}[q(h \mid x; \boldsymbol{\phi}) \| p(h; \boldsymbol{\theta}))] \quad (22)$$

where $p(h)$ is a Gaussian prior.

# Perplexity

| Model | Dim | 20News | RCV1 |
|-------|-----|--------|------|
| LDA | 50 | 1091 | 1437 |
| LDA | 200 | 1058 | 1142 |
| RSM | 50 | 953 | 988 |
| docNADE | 50 | 896 | 742 |
| SBN | 50 | 909 | 784 |
| fDARN | 50 | 917 | 724 |
| fDARN | 200 | —— | 598 |
| NVDM | 50 | **836** | 563 |
| NVDM | 200 | 852 | **550** |

[Miao et al., 2016]

# Topics

| Space | Religion | Encryption | Sport | Policy |
|-------|----------|------------|-------|--------|
| orbit | muslims | rsa | goals | bush |
| lunar | worship | cryptography | pts | resources |
| solar | belief | crypto | teams | charles |
| shuttle | genocide | keys | league | austin |
| moon | jews | pgp | team | bill |
| launch | islam | license | players | resolution |
| fuel | christianity | secure | nhl | mr |
| nasa | atheists | key | stats | misc |
| satellite | muslim | escrow | min | piece |
| japanese | religious | trust | buf | marc |

[Miao et al., 2016]

# Summary

1. Variational Inference, Review

2. Variational Auto-encoder

3. Application: Neural variational document model

# Reference

Kingma, D. P. and Welling, M. (2014).
Auto-encoding variational Bayes.
In *ICLR.*

Miao, Y., Yu, L., and Blunsom, P. (2016).
Neural variational inference for text processing.
In *International Conference on Machine Learning*, pages 1727–1736.