# Probabilistic Topic Models
## DGM4NLP

Miguel Rios
University of Amsterdam

May 5, 2019

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Outline

1 Probabilistic Topic Models

2 Neural Variational Inference for Text Processing

3 Discovering Discrete Latent Topics

4 LDA VAE

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Introduction

- Topic modelling provides models for automatically organizing, understanding, searching, and summarizing large corpus of documents.

[top]

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Introduction

- Topic modelling provides models for automatically organizing, understanding, searching, and summarizing large corpus of documents.
- Discover the hidden domains in the corpus.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
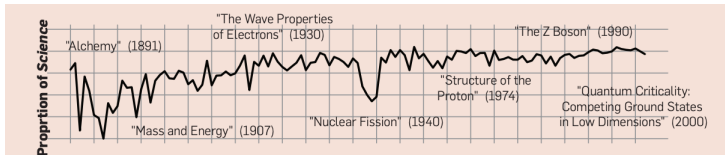LDA VAE
References
References

## Introduction

- Topic modelling provides models for automatically organizing, understanding, searching, and summarizing large corpus of documents.
- Discover the hidden domains in the corpus.
- Annotate the documents according to those domains.

[top]

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Introduction

- Topic modelling provides models for automatically organizing, understanding, searching, and summarizing large corpus of documents.
- Discover the hidden domains in the corpus.
- Annotate the documents according to those domains.
- Use annotations to organise, summarise, search, and make predictions over documents.
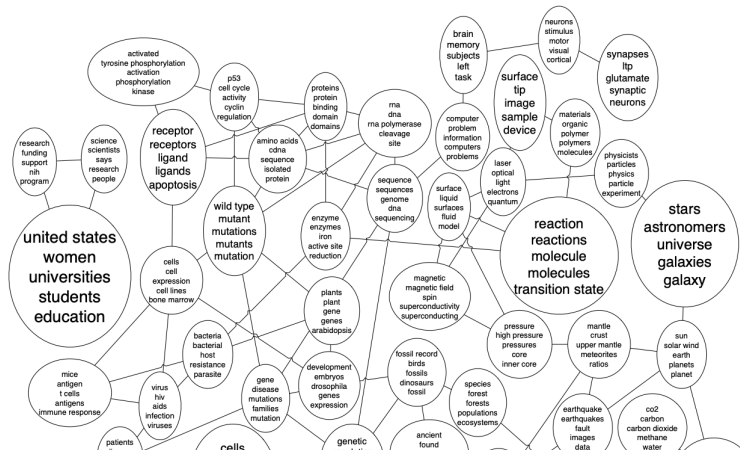
Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Probabilistic Topic Models

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Probabilistic Topic Models

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Probabilistic Topic Models

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Probabilistic Models

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Latent Dirichlet allocation (LDA)

- Motivation is that documents show multiple topics.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Latent Dirichlet allocation (LDA)

- Motivation is that documents show multiple topics.
- For example, in "Seeking Life's Bare (Genetic) Necessities," is about using data analysis to determine the number of genes an organism needs to survive (in an evolutionary sense).

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Latent Dirichlet allocation (LDA)

- Motivation is that documents show multiple topics.
- For example, in "Seeking Life's Bare (Genetic) Necessities," is about using data analysis to determine the number of genes an organism needs to survive (in an evolutionary sense).
- Highlighted words related to data analysis: **computer** and **prediction**, are highlighted in blue;
  and evolutionary biology: **life** and **organism**, in pink;

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Latent Dirichlet allocation (LDA)

- Motivation is that documents show multiple topics.
- For example, in "Seeking Life's Bare (Genetic) Necessities," is about using data analysis to determine the number of genes an organism needs to survive (in an evolutionary sense).
- Highlighted words related to data analysis: **computer** and **prediction**, are highlighted in blue;
  and evolutionary biology: **life** and **organism**, in pink;
- LDA is described by its generative process, the imaginary random process by which the model assumes the documents arose.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Latent Dirichlet allocation (LDA)

- Motivation is that documents show multiple topics.
- For example, in "Seeking Life's Bare (Genetic) Necessities," is about using data analysis to determine the number of genes an organism needs to survive (in an evolutionary sense).
- Highlighted words related to data analysis: **computer** and **prediction**, are highlighted in blue;
  and evolutionary biology: **life** and **organism**, in pink;
- LDA is described by its generative process, the imaginary random process by which the model assumes the documents arose.
- We denote a topic to be a distribution over a fixed vocabulary.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Latent Dirichlet allocation (LDA)

- Motivation is that documents show multiple topics.
- For example, in "Seeking Life's Bare (Genetic) Necessities," is about using data analysis to determine the number of genes an organism needs to survive (in an evolutionary sense).
- Highlighted words related to data analysis: **computer** and **prediction**, are highlighted in blue;
  and evolutionary biology: **life** and **organism**, in pink;
- LDA is described by its generative process, the imaginary random process by which the model assumes the documents arose.
- We denote a topic to be a distribution over a fixed vocabulary.
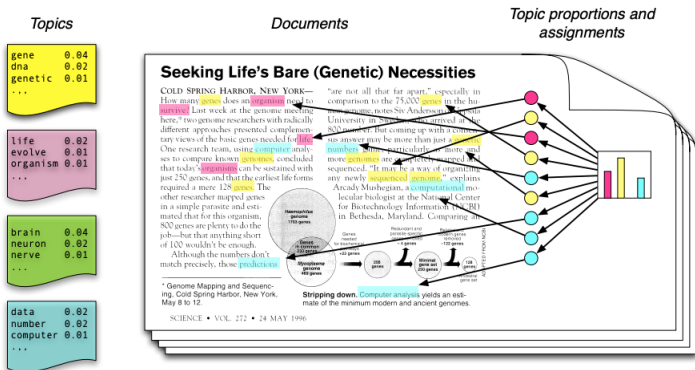- For example, the genetics topic contains words about genetics with high probability.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References
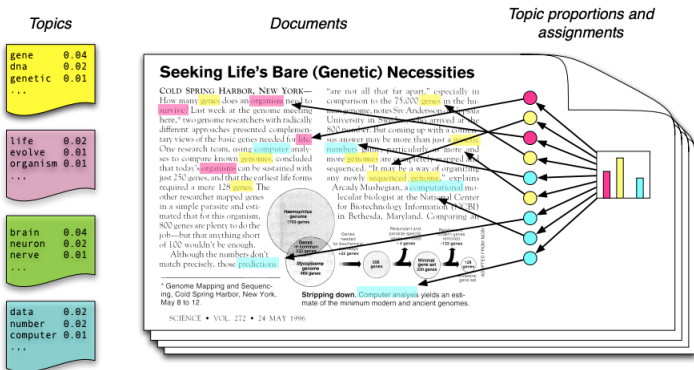
# Latent Dirichlet allocation (LDA)

- Motivation is that documents show multiple topics.
- For example, in "Seeking Life's Bare (Genetic) Necessities," is about using data analysis to determine the number of genes an organism needs to survive (in an evolutionary sense).
- Highlighted words related to data analysis: **computer** and **prediction**, are highlighted in blue; and evolutionary biology: **life** and **organism**, in pink;
- LDA is described by its generative process, the imaginary random process by which the model assumes the documents arose.
- We denote a topic to be a distribution over a fixed vocabulary.
- For example, the genetics topic contains words about genetics with high probability.
- We assume that these topics are specified before any data has

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# LDA



Topics

| gene | 0.04 |
| dna | 0.02 |
| genetic | 0.01 |
| ... | |

| life | 0.02 |
| evolve | 0.01 |
| organism | 0.01 |
| ... | |

| brain | 0.04 |
| neuron | 0.02 |
| nerve | 0.01 |
| ... | |

| data | 0.02 |
| number | 0.02 |
| computer | 0.01 |
| ... | |

Documents

Topic proportions and assignments

### Seeking Life's Bare (Genetic) Necessities

- Each topic is a distribution over words

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA



- Each topic is a distribution over words
- Each document is a mixture of corpus-wide topics

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# LDA



- Each topic is a distribution over words
- Each document is a mixture of corpus-wide topics

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA Objective



Topics     Documents     Topic proportions and assignments

- We only observe the documents

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA Objective



Topics · Documents · Topic proportions and assignments

- We only observe the documents
- The conditional distribution of the topic structure given the observed documents

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA

For each document :

1. Randomly choose a distribution over topics.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA

For each document :

1. Randomly choose a distribution over topics.
2. For each word in the document:

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA

For each document :

1. Randomly choose a distribution over topics.
2. For each word in the document:
   a Randomly choose a topic from the distribution over topics in step 1.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA

For each document :

1. Randomly choose a distribution over topics.
2. For each word in the document:
   a Randomly choose a topic from the distribution over topics in step 1.
   b Randomly choose a word from the corresponding distribution over the vocabulary

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
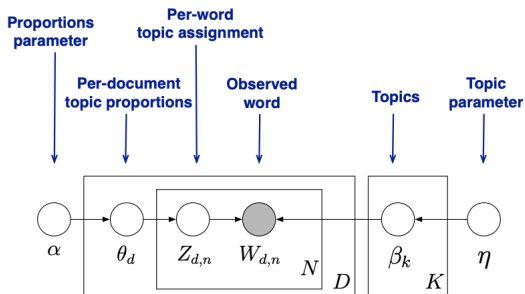LDA VAE
References
References

## LDA

For each document :

1. Randomly choose a distribution over topics.
2. For each word in the document:
   a Randomly choose a topic from the distribution over topics in step 1.
   b Randomly choose a word from the corresponding distribution over the vocabulary

- Each document exhibits the topics in different proportion (step1); each word in each document is drawn from one of the topics (step 2b), where the selected topic is chosen from the per-document distribution over topics (step 2a)

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA

For each document :

1. Randomly choose a distribution over topics.
2. For each word in the document:
   a. Randomly choose a topic from the distribution over topics in step 1.
   b. Randomly choose a word from the corresponding distribution over the vocabulary

- Each document exhibits the topics in different proportion (step1); each word in each document is drawn from one of the topics (step 2b), where the selected topic is chosen from the per-document distribution over topics (step 2a)
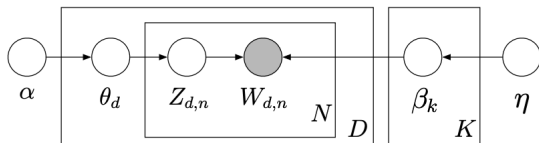- From the example article, the distribution over topics would place probability on genetics, data analysis, and evolutionary biology, and each word is drawn from one of those three

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# LDA PGM

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA PGM



- This joint defines a posterior, $p(\theta, z, \beta|w)$.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA PGM



- This joint defines a posterior, $p(\theta, z, \beta | w)$.
- From a collection of documents, infer

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA PGM



- This joint defines a posterior, $p(\theta, z, \beta | w)$.
- From a collection of documents, infer
- Per-word topic assignment $z_{d,n}$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA PGM



- This joint defines a posterior, $p(\theta, z, \beta | w)$.
- From a collection of documents, infer
- Per-word topic assignment $z_{d,n}$
- Per-document topic proportions $\theta_d$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA PGM



- This joint defines a posterior, $p(\theta, z, \beta | w)$.
- From a collection of documents, infer
- Per-word topic assignment $z_{d,n}$
- Per-document topic proportions $\theta_d$
- Per-corpus topic distributions $\beta_k$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
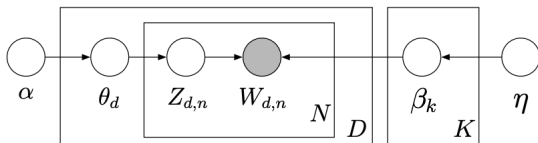Discovering Discrete Latent Topics
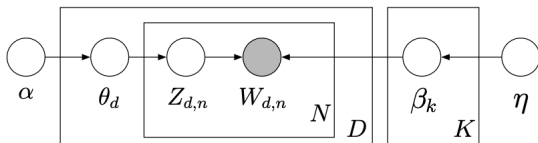LDA VAE
References
References

## LDA PGM



- This joint defines a posterior, $p(\theta, z, \beta | w)$.
- From a collection of documents, infer
- Per-word topic assignment $z_{d,n}$
- Per-document topic proportions $\theta_d$
- Per-corpus topic distributions $\beta_k$
- Then use posterior expectations to perform the task at hand: information retrieval, document similarity, exploration, and others

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Dirichlet distribution

- The Dirichlet distribution is an exponential family distribution over the simplex, i.e., positive vectors that sum to one

$$p(\theta|\alpha) = \frac{\Gamma\left(\sum_i \alpha_i\right)}{\prod_i \Gamma\left(\alpha_i\right)} \prod_i \theta_i^{\alpha_i - 1} \tag{1}$$

- It is conjugate to the multinomial. Given a multinomial observation, the posterior distribution of $\theta$ is a Dirichlet.
- The parameter $\alpha$ controls the mean shape and sparsity of $\theta$.
- The topic proportions are a K dimensional Dirichlet. The topics are a V dimensional Dirichlet.
- The alpha controls the mixture of topics for any given document.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# Dirichlet distribution

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Dirichlet distribution

- At low alpha values (less than one), most of the topic distribution samples are in the corners (near the topics).

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Dirichlet distribution

- At low alpha values (less than one), most of the topic distribution samples are in the corners (near the topics).

- At alpha equal to one, any space on the surface of the triangle (3-simplex) is fair game (uniformly distributed). You could equally likely end up with a sample favoring only one topic, a sample that gives an even mixture of all the topics, or something in between.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Dirichlet distribution

- At low alpha values (less than one), most of the topic distribution samples are in the corners (near the topics).

- At alpha equal to one, any space on the surface of the triangle (3-simplex) is fair game (uniformly distributed). You could equally likely end up with a sample favoring only one topic, a sample that gives an even mixture of all the topics, or something in between.
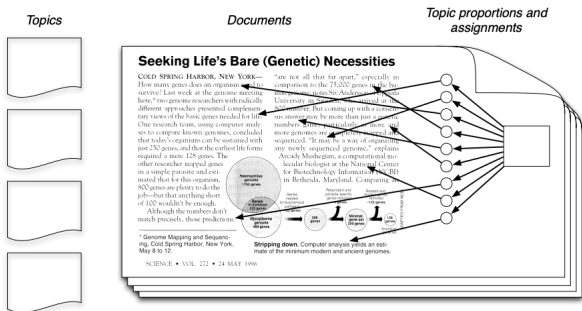
- For alpha values greater than one, the samples start to congregate to the center. This means that as alpha gets bigger, your samples will more likely be uniform or an even mixture of all the topics.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Working

LDA trades off two goals.

- **(1)** For each document, allocate its words to as few topics as possible.

  **(2)** For each topic, assign high probability to as few terms as possible.

- Putting a document in a single topic makes 2 hard: All of its words must have probability under that topic.

- Putting very few words in each topic makes 1 hard: To cover a document's words, it must assign many topics to it.

- Trading off these goals finds groups of tightly co-occurring words.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Posterior Inference



Topics       Documents       Topic proportions and assignments

- Our goal is to compute the distribution of the hidden variables conditioned on the documents
  p(topics, proportions, assignments—documents)

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Posterior Inference



- The joint distribution of the latent variables and documents is
  $\prod_{i=1}^{K} p(\beta_i|\eta) \prod_{d=1}^{D} p(\theta_d|\alpha) \left( \prod_{n=1}^{N} p(z_{d,n}|\theta_d) p(w_{\alpha,n}|\beta_{1:k,z_{d,n}}) \right)$
- The posterior of the latent variables given the documents is
  $p(\beta, \theta, \mathbf{z}|\mathbf{w})$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Posterior Inference

- $p(\beta, \theta, \mathbf{z}|\mathbf{w}) = \frac{p(\beta, \theta, \mathbf{z}, \mathbf{w})}{\int_\beta \int_\theta \sum_{\mathbf{z}} p(\beta, \theta, \mathbf{z}, \mathbf{w})}$

- The denominator, the marginal $p(w)$ is intractable

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Posterior Inference



*Sample one document*   *Analyze it*   *Update the model*

- Condition on large data sets and approximate the posterior.
- Variational inference, we optimize over a family of distributions to find the member closest in KL divergence to the posterior.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Posterior Inference



*Sample one document*   *Analyze it*   *Update the model*

1. Sample a document $w_d$ from the collection
2. Infer how $w_d$ exhibits the current topics
3. Create intermediate topics, formed as though the $w_d$ is the only document.
4. Adjust the current topics according to the intermediate topics.
5. Repeat.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Mean-field variational inference for LDA



1. Document variables: Topic proportions $\theta$ and topic assignments $z_{1:N}$ .

2. Corpus variables: Topics $\beta_{1:K}$

3. The variational approximation is:
$q(\boldsymbol{\beta}, \boldsymbol{\theta}, z) =$
$\prod_{k=1}^{K} q\left(\beta_k | \lambda_k\right) \prod_{d=1}^{D} q\left(\theta_d | \gamma_d\right) \prod_{n=1}^{N} q\left(z_{d,n} | \phi_{d,n}\right)$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Mean-field variational inference for LDA

1: Initialize topics randomly.
2: **repeat**
3:    **for** each document **do**
4:       **repeat**
5:          Update the topic assignment variational parameters.
6:          Update the topic proportions variational parameters.
7:       **until** document objective converges
8:    **end for**
9:    Update the topics from aggregated per-document parameters.
10: **until** corpus objective converges.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# LDA Extensions

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Correlated topic models



- Draw topic proportions from a logistic normal

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Correlated topic models



- Draw topic proportions from a logistic normal
- Allows topic occurrences to have correlation.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Correlated topic models



- Draw topic proportions from a logistic normal
- Allows topic occurrences to have correlation.
- Gives a map of topics and how they are related

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Correlated topic models



- Draw topic proportions from a logistic normal
- Allows topic occurrences to have correlation.
- Gives a map of topics and how they are related
- Better fit for observed data, but computation is more complex

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Dynamic topic models

- LDA assumes that the order of documents does not matter.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Dynamic topic models

- LDA assumes that the order of documents does not matter.
- Corpora span hundreds of years

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Dynamic topic models

- Each document has an influence score *ld*.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Dynamic topic models

- Each document has an influence score $ld$.
- Each topic is biased with the documents with high influence.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Dynamic topic models

- Each document has an influence score $Id$.
- Each topic is biased with the documents with high influence.
- The posterior of the influence scores could find articles that best explain the changes in language.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# Dynamic topic models

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Outline

1. Probabilistic Topic Models

2. Neural Variational Inference for Text Processing

3. Discovering Discrete Latent Topics

4. LDA VAE

Probabilistic Topic Models
**Neural Variational Inference for Text Processing**
Discovering Discrete Latent Topics
LDA VAE
References
References

## Neural Variational Inference for Text Processing

- Neural variational framework for generative models of documents based on the variational auto-encoder.

[Miao et al., 2016]

Probabilistic Topic Models
**Neural Variational Inference for Text Processing**
Discovering Discrete Latent Topics
LDA VAE
References
References

## Neural Variational Inference for Text Processing

- Neural variational framework for generative models of documents based on the variational auto-encoder.
- NVDM is a generative model of text which aims to extract a continuous semantic latent variable for each document.

[Miao et al., 2016]

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# Neural Variational Inference for Text Processing

- Neural variational framework for generative models of documents based on the variational auto-encoder.
- NVDM is a generative model of text which aims to extract a continuous semantic latent variable for each document.
- Model is denoted by variational auto-encoder:

[Miao et al., 2016]

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Neural Variational Inference for Text Processing

- Neural variational framework for generative models of documents based on the variational auto-encoder.
- NVDM is a generative model of text which aims to extract a continuous semantic latent variable for each document.
- Model is denoted by variational auto-encoder:
- MLP encoder (inference) compresses the bag-of-words document representation into a continuous latent distribution,

[Miao et al., 2016]

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Neural Variational Inference for Text Processing

- Neural variational framework for generative models of documents based on the variational auto-encoder.
- NVDM is a generative model of text which aims to extract a continuous semantic latent variable for each document.
- Model is denoted by variational auto-encoder:
- MLP encoder (inference) compresses the bag-of-words document representation into a continuous latent distribution,
- Softmax decoder (generative model) reconstructs the document by generating the words independently.

[Miao et al., 2016]

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# Neural Variational Inference for Text Processing

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# Neural Variational Inference for Text Processing

- Let $X in R^{|V|}$ be the bag-of-words representation of a document and $x_i in R^{|V|}$ be the one-hot representation of the word at position $i$.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# Neural Variational Inference for Text Processing

- Let $X in R^{|V|}$ be the bag-of-words representation of a document and $x_i in R^{|V|}$ be the one-hot representation of the word at position $i$.
- MLP encoder $q(z|x)$ compresses document representations into continuous hidden vectors

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# Neural Variational Inference for Text Processing

- Let $X in R^{|V|}$ be the bag-of-words representation of a document and $x_i in R^{|V|}$ be the one-hot representation of the word at position $i$.
- MLP encoder $q(z|x)$ compresses document representations into continuous hidden vectors
- Softmax decoder $p(\mathbf{x}|\mathbf{z}) = \prod_{i=1}^{N} p(\mathbf{x}_i|\mathbf{z})$ reconstructs the documents by independently generating the words.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Neural Variational Inference for Text Processing

- Let $X in R^{|V|}$ be the bag-of-words representation of a document and $x_i in R^{|V|}$ be the one-hot representation of the word at position $i$.
- MLP encoder $q(z|x)$ compresses document representations into continuous hidden vectors
- Softmax decoder $p(x|z) = \prod_{i=1}^{N} p(x_i|z)$ reconstructs the documents by independently generating the words.
- We derive the lower bound:
  $\mathcal{L} = \mathbb{E}_{q_\phi(z|x)} \left[ \sum_{i=1}^{N} \log p_\theta(x_i|z) \right] - D_{\mathrm{KL}}(q_\phi(z|x) \| p(z))$
  where N is the number of words in the document

Probabilistic Topic Models
**Neural Variational Inference for Text Processing**
Discovering Discrete Latent Topics
LDA VAE
References
References

## Data

- Standard news corpora:

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Data

- Standard news corpora:

- 20NewsGroups is a collection of newsgroup documents, consisting of 11,314 training and 7,531 test articles.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Data

- Standard news corpora:
- 20NewsGroups is a collection of newsgroup documents, consisting of 11,314 training and 7,531 test articles.
- Reuters RCV1-v2 is a large collection from Reuters newswire stories with 794,414 training and 10,000 test cases.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Data

- Standard news corpora:
- 20NewsGroups is a collection of newsgroup documents, consisting of 11,314 training and 7,531 test articles.
- Reuters RCV1-v2 is a large collection from Reuters newswire stories with 794,414 training and 10,000 test cases.
- The vocabulary size of these two datasets are set as 2,000 and 10,000

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Results

| Model | Dim | 20News | RCV1 |
|-------|-----|--------|------|
| LDA   | 50  | 1091   | 1437 |
| LDA   | 200 | 1058   | 1142 |
| NVDM  | 50  | **836** | 563 |
| NVDM  | 200 | 852    | **550** |

- perplexity is computed $ppl = \exp\left(-\frac{1}{D}\sum_n^{N_d}\frac{1}{N_d}\log p\left(\mathbf{x}_d\right)\right)$, where D is the number of documents, $N_d$ represents the length of the dth document.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Results

| Model | Dim | 20News | RCV1 |
|-------|-----|--------|------|
| LDA   | 50  | 1091   | 1437 |
| LDA   | 200 | 1058   | 1142 |
| NVDM  | 50  | **836**  | 563  |
| NVDM  | 200 | 852    | **550** |

- perplexity is computed $ppl = \exp\left(-\frac{1}{D}\sum_{n}^{N_d}\frac{1}{N_d}\log p\left(\mathbf{x}_d\right)\right)$, where D is the number of documents, $N_d$ represents the length of the dth document.
- Since $logp(x)$ in the NVDM is the variational lower bound (which is an upper bound on perplexity).

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Results

The topics learned by NVDM on 20News

| *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 |

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Outline

1. Probabilistic Topic Models

2. Neural Variational Inference for Text Processing

3. Discovering Discrete Latent Topics

4. LDA VAE

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics

- Introduce a neural network to parameterise the multinomial topic distribution

$$\theta_d \sim G\left(\mu_0, \sigma_0^2\right), \text{ for } d \in D$$
$$z_n \sim \text{Multi}\left(\theta_d\right), \text{ for } n \in [1, N_d] \qquad (2)$$
$$w_n \sim \text{Multi}\left(\beta_{z_n}\right), \text{ for } n \in [1, N_d]$$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics

- Introduce a neural network to parameterise the multinomial topic distribution

$$\theta_d \sim \mathrm{G}\left(\mu_0, \sigma_0^2\right), \text{ for } d \in D$$
$$z_n \sim \mathrm{Multi}\left(\theta_d\right), \text{ for } n \in [1, N_d] \quad (2)$$
$$w_n \sim \mathrm{Multi}\left(\beta_{z_n}\right), \text{ for } n \in [1, N_d]$$

- $G(\mu_0, \sigma_0^2)$ is composed of a NN $\theta = g(x)$ conditioned on a isotropic Gaussian $x \sim \mathcal{N}(\mu_0, \sigma_0^2 0)$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics

- Introduce a neural network to parameterise the multinomial topic distribution

$$\theta_d \sim \mathrm{G}\left(\mu_0, \sigma_0^2\right), \text{ for } d \in D$$
$$z_n \sim \text{Multi}\left(\theta_d\right), \text{ for } n \in [1, N_d] \quad (2)$$
$$w_n \sim \text{Multi}\left(\beta_{z_n}\right), \text{ for } n \in [1, N_d]$$

- $G(\mu_0, \sigma_0^2)$ is composed of a NN $\theta = g(x)$ conditioned on a isotropic Gaussian $x \sim \mathcal{N}(\mu_0, \sigma_0^2 0)$

- Gaussian Softmax Construction pass a Gaussian random vector through a softmax function to parameterise the multinomial document topic distributions.

$$x \sim \mathcal{N}\left(\mu_0, \sigma_0^2\right)$$
$$\theta = \text{softmax}\left(W_1^T x\right) \quad (3)$$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# Discovering Discrete Latent Topics



[Miao et al., 2017]

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics

- Neural Topic Models with a finite number of topics K.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics

- Neural Topic Models with a finite number of topics K.
- The topic distribution over words given a topic assignment $z_n$ is
  $p(w_n|\beta, z_n) = \text{Multi}(\beta_{z_n}).$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics

- Neural Topic Models with a finite number of topics K.
- The topic distribution over words given a topic assignment $z_n$ is
  $p(w_n | \beta, z_n) = \text{Multi}(\beta_{z_n})$.
- Introduce topic vectors $t \in R^{K \times H}$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics

- Neural Topic Models with a finite number of topics K.
- The topic distribution over words given a topic assignment $z_n$ is
  $p(w_n|\beta, z_n) = \text{Multi}(\beta_{z_n})$.
- Introduce topic vectors $t \in R^{K \times H}$
- word vectors $v \in R^{V \times H}$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics

- Neural Topic Models with a finite number of topics K.
- The topic distribution over words given a topic assignment $z_n$ is
  $p(w_n|\beta, z_n) = \text{Multi}(\beta_{z_n})$.
- Introduce topic vectors $t \in R^{K \times H}$
- word vectors $v \in R^{V \times H}$
- and generate the topic distributions over words by:
  $\beta_k = \text{softmax}\left(v \cdot t_k^T\right)$
  $\beta \in R^{K \times V}$ is the semantic similarity between topics and words.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics

- With lower bound:
$\mathcal{L}_d = \sum_{n=1}^{N} \left[\log p\left(w_n|\beta,\hat{\theta}\right)\right] - D_{KL}[q(x|d)\|p(x)]$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# Discovering Discrete Latent Topics

- With lower bound:
  $\mathcal{L}_d = \sum_{n=1}^{N} \left[ \log p\left(w_n | \beta, \hat{\theta}\right) \right] - D_{KL}[q(x|d) \| p(x)]$

- 
$$
\begin{aligned}
\log p\left(w_n | \beta, \hat{\theta}\right) &= \log \sum_{z_n} \left[ p\left(w_n | \beta_{z_n}\right) p\left(z_n | \hat{\theta}\right) \right] \\
&= \log(\hat{\theta} \cdot \beta)
\end{aligned}
\tag{4}
$$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics

- With lower bound:
  $\mathcal{L}_d = \sum_{n=1}^{N} \left[ \log p \left( w_n | \beta, \hat{\theta} \right) \right] - D_{KL}[q(x|d) \| p(x)]$

-
$$\log p \left( w_n | \beta, \hat{\theta} \right) = \log \sum_{z_n} \left[ p \left( w_n | \beta_{z_n} \right) p \left( z_n | \hat{\theta} \right) \right]$$
$$= \log(\hat{\theta} \cdot \beta) \qquad (4)$$

CON addition of topic diversity regularisation to the objective

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics

- Unbounded neural topic models the topics $t \in R^{\infty \times H}$ are dynamically generated by $RNN_{Topic}$ The generation of $\beta$ is as follows:

$$
\begin{aligned}
t_k &= RNN_{Topic}\ (t_{k-1}) \\
\beta_k &= \text{softmax}\left(v \cdot t_k^T\right)
\end{aligned}
\tag{5}
$$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics

- Unbounded neural topic models the topics $t \in R^{\infty \times H}$ are dynamically generated by $\text{RNN}_{\text{Topic}}$ The generation of $\beta$ is as follows:

$$
\begin{aligned}
t_k &= \text{RNN }_{\text{Topic }} (t_{k-1}) \\
\beta_k &= \text{softmax}\left(v \cdot t_k^T\right)
\end{aligned}
\tag{5}
$$

- where $v$ represents the word vectors, $t_k$ is the kth topic generated by RNN

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
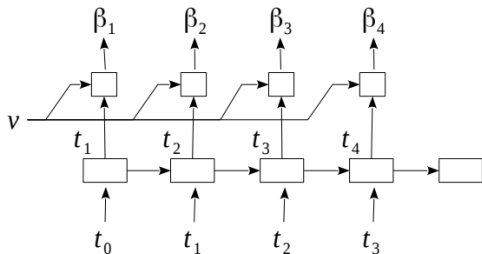References
References

## Discovering Discrete Latent Topics

- Unbounded neural topic models the topics $t \in R^{\infty \times H}$ are dynamically generated by $\text{RNN}_{\text{Topic}}$ The generation of $\beta$ is as follows:

$$t_k = \text{RNN}_{\text{Topic}} (t_{k-1})$$
$$\beta_k = \text{softmax} \left( v \cdot t_k^T \right) \tag{5}$$

- where $v$ represents the word vectors, $t_k$ is the kth topic generated by RNN

- If $I > \gamma$, we increase the active number of topics $i$ by 1, $\mathcal{I} = \sum_d^D \left[ \mathcal{L}_d^i - \mathcal{L}_d^{i-1} \right] / \sum_d^D \left[ \mathcal{L}_d^i \right]$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Discovering Discrete Latent Topics



- 

$$t_k = \text{RNN}_{\text{Topic}}(t_{k-1})$$
$$\beta_k = \text{softmax}\left(v \cdot t_k^T\right)$$

(6)

Probabilistic Topic Models
Neural Variational Inference for Text Processing
**Discovering Discrete Latent Topics**
LDA VAE
References
References

## Results

| **Finite Topic Model** | MXM | | 20News | | RCV1 | |
|---|---|---|---|---|---|---|
| | 50 | 200 | 50 | 200 | 50 | 200 |
| GSM | **306** | **272** | **822** | 830 | **717** | **602** |
| GSB | 309 | 296 | 838 | 826 | 788 | 634 |
| RSB | 311 | 297 | 835 | **822** | 750 | 628 |
| OnlineLDA | 312 | 342 | 893 | 1015 | 1062 | 1058 |
| (Hoffman et al., 2010) | | | | | | |
| NVLDA | 330 | 357 | 1073 | 993 | 791 | 797 |
| (Srivastava & Sutton, 2016) | | | | | | |

| **Unbounded Topic Model** | MXM | 20News | RCV1 |
|---|---|---|---|
| RSB-TF | **303** | **825** | **622** |
| HDP (Wang et al., 2011) | 370 | 937 | 918 |

- *MXM the Million Song Dataset with 210,519 training and

Probabilistic Topic Models
Neural Variational Inference for Text Processing
**Discovering Discrete Latent Topics**
LDA VAE
References
References

## Results

| *Space* | *Religion* | *Encryption* | *Sport* | *Science* |
|---------|-----------|--------------|---------|-----------|
| space | god | encryption | player | science |
| satellite | atheism | device | hall | theory |
| april | exist | technology | defensive | scientific |
| sequence | atheist | protect | team | universe |
| launch | moral | americans | average | experiment |
| president | existence | chip | career | observation |
| station | marriage | use | league | evidence |
| radar | system | privacy | play | exist |
| training | parent | industry | bob | god |
| committee | murder | enforcement | year | mistake |

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# Outline

1. Probabilistic Topic Models

2. Neural Variational Inference for Text Processing

3. Discovering Discrete Latent Topics

4. LDA VAE

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA VAE

- Effective VAE based model for LDA

[Srivastava and Sutton, 2017]

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA VAE

- Effective VAE based model for LDA
- Dirichlet within VAE is difficult to develop an effective reparameterisation function
  Solve by constructing a Laplace approximation to the Dirichlet prior.

[Srivastava and Sutton, 2017]

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA VAE

- Effective VAE based model for LDA
- Dirichlet within VAE is difficult to develop an effective reparameterisation function
  Solve by constructing a Laplace approximation to the Dirichlet prior.
- This approximation to the Dirichlet results in the distribution over the softmax variables

[Srivastava and Sutton, 2017]

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Laplace approximation

- Approximation in the softmax basis instead of the simplex.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Laplace approximation

- Approximation in the softmax basis instead of the simplex.
- Dirichlet probability density function over the softmax variable $h$ is:

$$P(\theta(\mathbf{h})|\alpha) = \frac{\Gamma\left(\sum_k \alpha_k\right)}{\prod_k \Gamma\left(\alpha_k\right)} \prod_k \theta_k^{\alpha_k} g\left(\mathbf{1}^T \mathbf{h}\right) \qquad (7)$$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Laplace approximation

- Approximation in the softmax basis instead of the simplex.
- Dirichlet probability density function over the softmax variable $h$ is:

$$P(\theta(\mathbf{h})|\alpha) = \frac{\Gamma\left(\sum_k \alpha_k\right)}{\prod_k \Gamma\left(\alpha_k\right)} \prod_k \theta_k^{\alpha_k} g\left(\mathbf{1}^T \mathbf{h}\right) \tag{7}$$

- Here $\theta = \sigma(h)$, where $\sigma(\cdot)$ represents the softmax function

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA VAE

- Approximation to the Dirichlet results in the distribution over the softmax variables h as a multivariate normal with mean $\mu_1$ and covariance matrix $\Sigma_1$ where:

$$
\begin{aligned}
\mu_{1k} &= \log \alpha_k - \frac{1}{K} \sum_i \log \alpha_i \\
\Sigma_{1kk} &= \frac{1}{\alpha_k} \left( 1 - \frac{2}{K} \right) + \frac{1}{K^2} \sum_i \frac{1}{\alpha_k}
\end{aligned}
\tag{8}
$$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA VAE

- Approximate of $p(\theta|\alpha)$ with $\hat{p}(\theta|\mu_1, \Sigma_1) = \mathcal{LN}(\theta|\mu_1, \Sigma_1)$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA VAE

- Approximate of $p(\theta|\alpha)$ with $\hat{p}(\theta|\mu_1, \Sigma_1) = \mathcal{LN}(\theta|\mu_1, \Sigma_1)$
- where LN is a logistic normal distribution with parameters $\mu_1$, $\Sigma_1$ for k (number of topics).

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## LDA VAE

- Approximate of $p(\theta|\alpha)$ with $\hat{p}(\theta|\mu_1, \Sigma_1) = \mathcal{LN}(\theta|\mu_1, \Sigma_1)$
- where LN is a logistic normal distribution with parameters $\mu_1$, $\Sigma_1$ for k (number of topics).
- and ELBO:

$$L(\Theta) = \sum_{d=1}^{D} \left[ -\left( \frac{1}{2} \left\{ \text{tr}\left(\Sigma_1^{-1}\Sigma_0\right) + \left(\mu_1 - \mu_0\right)^T \Sigma_1^{-1} \left(\mu_1 - \mu_0\right) - K + \log \frac{|\Sigma_1|}{|\Sigma_0|} \right\} \right. \right.$$
$$\left. \left. + \mathbb{E}_{\epsilon \sim \mathcal{N}(0,I)} \left[ \mathbf{w}_d^\top \log \left( \sigma(\boldsymbol{\beta})\sigma\left(\mu_0 + \Sigma_0^{1/2}\epsilon\right) \right) \right. \right. \right.$$
$$(9)$$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
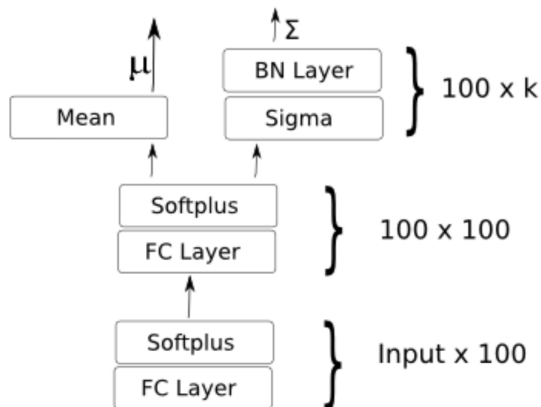LDA VAE
References
References

## LDA VAE

- Approximate of $p(\theta|\alpha)$ with $\hat{p}(\theta|\mu_1, \Sigma_1) = \mathcal{LN}(\theta|\mu_1, \Sigma_1)$
- where LN is a logistic normal distribution with parameters $\mu_1$, $\Sigma_1$ for k (number of topics).
- and ELBO:

$$L(\Theta) = \sum_{d=1}^{D}\left[-\left(\frac{1}{2}\left\{\text{tr}\left(\Sigma_1^{-1}\Sigma_0\right) + \left(\mu_1 - \mu_0\right)^T\Sigma_1^{-1}\left(\mu_1 - \mu_0\right) - K + \log\frac{|\Sigma_1|}{|\Sigma_0|}\right\}\right.\right.$$
$$\left.\left.+\mathbb{E}_{\epsilon\sim\mathcal{N}(0,I)}\left[\mathbf{w}_d^\top\log\left(\sigma(\boldsymbol{\beta})\sigma\left(\mu_0 + \Sigma_0^{1/2}\epsilon\right)\right)\right.\right.\right.$$
(9)

- with $\mu_0 = f_\mu(\mathbf{w}, \boldsymbol{\delta})$ and $\Sigma_0 = \text{diag}\left(f_\Sigma(\mathbf{w}, \boldsymbol{\delta})\right)$

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Architecture

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

## Results

| # topics | ProdLDA VAE | LDA VAE | LDA DMFVI | LDA Collapsed Gibbs | NVDM |
|----------|-------------|---------|-----------|---------------------|------|
| **50**   | 1172        | 1059    | 1046      | **728**             | 837  |
| **200**  | 1168        | 1128    | 1195      | **688**             | 884  |

ppl 20 Newsgroups

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
References

# Results

| Model | Topics |
|---|---|
| **ProdLDA** | motherboard meg printer quadra hd windows processor vga mhz connector<br>armenian genocide turks turkish muslim massacre turkey armenians armenia greek<br>voltage nec outlet circuit cable wiring wire panel motor install<br>season nhl team hockey playoff puck league flyers defensive player<br>israel israeli lebanese arab lebanon arabs civilian territory palestinian militia |
| **LDA NVLDA** | db file output program line entry write bit int return<br>drive disk get card scsi use hard ide controller one<br>game team play win year player get think good make<br>use law state health file gun public issue control firearm<br>people say one think life make know god man see |
| **LDA DMFVI** | write article dod ride right go get night dealer like<br>gun law use drug crime government court criminal firearm control<br>lunar flyers hitter spacecraft power us existence god go mean<br>stephanopoulos encrypt spacecraft ripem rsa cipher saturn violate lunar crypto<br>file program available server version include software entry ftp use |
| **LDA Collapsed Gibbs** | get right back light side like see take time one<br>list mail send post anonymous internet file information user message<br>thanks please know anyone help look appreciate get need email<br>jesus church god law say christian one christ day come<br>bike dod ride dog motorcycle write article bmw helmet get |
| **NVDM** | light die burn body life inside mother tear kill christian<br>insurance drug different sport friend bank owner vancouver buy prayer<br>input package interface output tape offer component channel level model<br>price quadra hockey slot san playoff jose deal market dealer<br>christian church gateway catholic christianity homosexual resurrection modem mouse sunday |

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
**References**

## References I

topic-models. URL http://www.cs.columbia.edu/~blei/papers/BleiLafferty2009.pdf.

Yishu Miao, Lei Yu, and Phil Blunsom. Neural variational inference for text processing. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1727–1736, New York, New York, USA, 20–22 Jun 2016. PMLR. URL http://proceedings.mlr.press/v48/miao16.html.

Probabilistic Topic Models
Neural Variational Inference for Text Processing
Discovering Discrete Latent Topics
LDA VAE
References
**References**

## References II

Yishu Miao, Edward Grefenstette, and Phil Blunsom. Discovering discrete latent topics with neural variational inference. *CoRR*, abs/1706.00359, 2017. URL http://arxiv.org/abs/1706.00359.

Akash Srivastava and Charles A. Sutton. Autoencoding variational inference for topic models. In *ICLR*, 2017.