

# CS 6501 Natural Language Processing

## Latent Variable Models

---

Yangfeng Ji

November 5, 2018

Department of Computer Science  
University of Virginia



ENGINEERING

# Overview

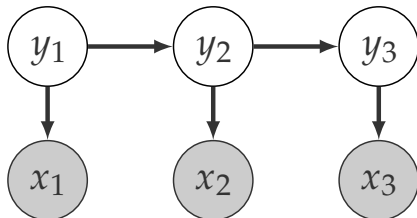
1. Latent Variable Models
2. Variational Inference
3. Example: Latent Dirichlet Allocation

# Latent Variable Models

---

# Latent Variables Models

## Hidden Markov Models



# Gaussian Mixture Models

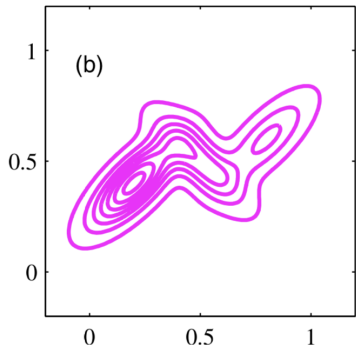
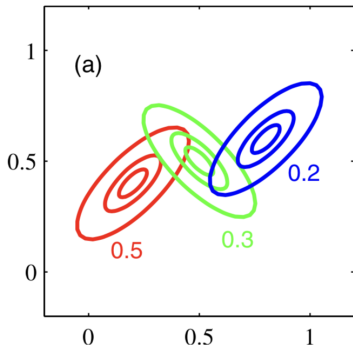
a Gaussian mixture model with  $K$  components and each component is a Gaussian distribution

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

Parameters

- ▶  $\boldsymbol{\mu}_k$ : mean of the  $k$ -th component
- ▶  $\boldsymbol{\Sigma}_k$ : variance of the  $k$ -th component
- ▶  $\pi_k$ : weight of the  $k$ -th component with  $\sum_k \pi_k = 1$

# Gaussian Mixture Models: Example



[Bishop, 2006]

# GMM as a Latent Variable Model

Define a  $K$ -dimensional binary random vector  $\mathbf{z}$  to indicate which mixture component a data point comes from

- ▶ only one component of  $\mathbf{z}$  is 1 and all the rest are 0.
- ▶ the probability of  $z_k$  is defined as

$$p(z_k = 1) = \pi_k \quad (2)$$

For each  $z_k$

$$p(z_k = 1) = \pi_k \quad (3)$$

Overall,

$$p(z) = \prod_{k=1}^K \pi_k^{z_k} \quad (4)$$

is a **categorical** distribution with parameters  $\{\pi_k\}$



Using  $\mathbf{z}$  as an indicator vector, we can redefine  $p(\mathbf{x} \mid \mathbf{z})$

►  $p(\mathbf{x} \mid z_k = 1)$

$$p(\mathbf{x} \mid z_k = 1) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (5)$$

Using  $\mathbf{z}$  as an indicator vector, we can redefine  $p(\mathbf{x} \mid \mathbf{z})$

►  $p(\mathbf{x} \mid z_k = 1)$

$$p(\mathbf{x} \mid z_k = 1) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (5)$$

►  $p(\mathbf{x} \mid \mathbf{z})$

$$p(\mathbf{x} \mid \mathbf{z}) = \sum_{k=1}^K \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \quad (6)$$

# Joint Probability

$$\begin{aligned} p(\mathbf{x}, \mathbf{z}) &= p(\mathbf{z})p(\mathbf{x} \mid \mathbf{z}) \\ &= \prod_{k=1}^K \pi_k^{z_k} \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \end{aligned} \tag{7}$$

Marginal

$$\begin{aligned} p(\mathbf{x}) &= \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) \\ &= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{aligned} \tag{8}$$

# Graphical Representation

With  $N$  data points from the GMM

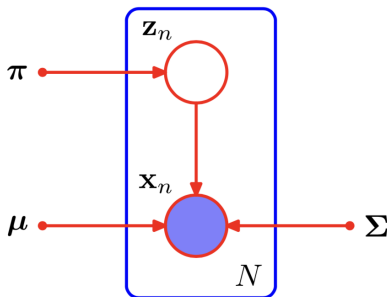
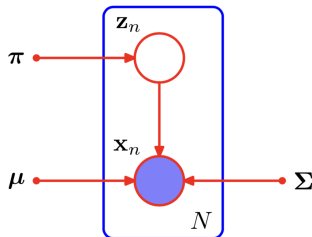


Figure: Graphical representation of GMM [Bishop, 2006].

# Generative Story



The generative story of a GMM can be formulated as

1. Randomly pick a mixture component  $k$ , with  $p(z_k = 1) = \pi_k$
2. Randomly generate a data point from the  $k$  component,  $\mathcal{N}(\mu_k, \Sigma_k)$

The procedure can be repeated multiple times

# Parameter Estimation

Given  $N$  data points  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , parameter estimation on a GMM is an iteration between the following two steps

1. Estimate  $p(\mathbf{z}_n)$  for every  $\mathbf{x}_n$
2. Estimate  $\{(\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\}_{k=1}^K$  based on  $\{p(\mathbf{z}_n)\}$  and  $\{\mathbf{x}_n\}$

Go back to step 1, until convergence

# Example

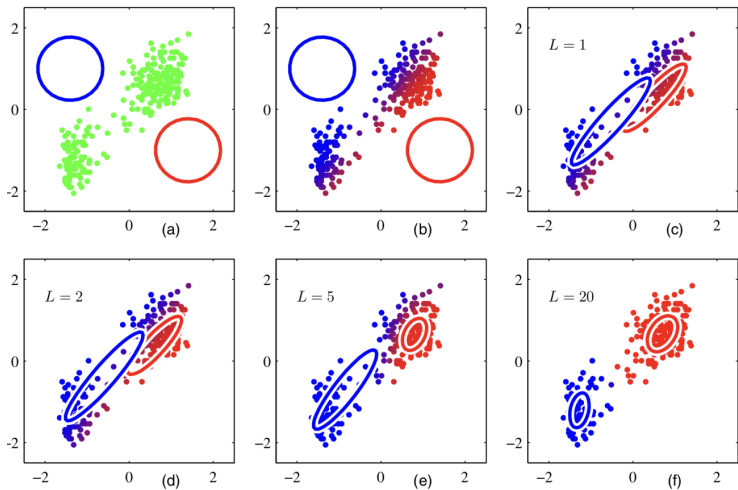
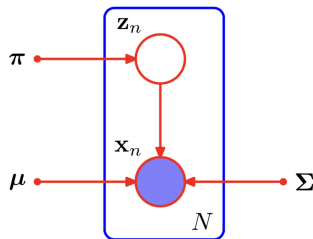


Figure: Illustration of the EM algorithm for GMM parameter estimation.

# Comments



A few things about latent variable formulation of GMMs

- ▶  $z_n$  is defined on each data point
- ▶ the model can be interpreted as a generative story
- ▶ marginalizing over  $z$  in  $p(x, z)$  is tractable



# Variational Inference

---

# Ideal Cases

We can solve a latent variable model by compute the following posterior distribution

$$p(\mathbf{z} \mid \mathbf{x}; \boldsymbol{\theta}) = \frac{p(\mathbf{z}, \mathbf{x})}{p(\mathbf{x})} \quad (9)$$

For example, MLE

$$\arg \max_{\boldsymbol{\theta}} \log p(\mathbf{z} \mid \mathbf{x}; \boldsymbol{\theta}) \quad (10)$$

However, the challenge comes from

$$p(x) = \sum_z p(x, z) \quad (11)$$

- ▶ Integral may be intractable, when  $z$  is continuous
- ▶ The space of  $z$  could be exponentially large, when  $z$  is discrete

BTW, in Bayesian statistics,  $p(x)$  is called **evidence**, it measure how good it is for a given model (parameter)

# Variational Inference

Instead of computing  $p(z \mid x)$ , we define a family of distribution  $\mathbb{Q}$ , and compute the following optimization problem

$$\tilde{q}(z) = \arg \min_{q(z) \in \mathbb{Q}} \text{KL}(q(z) \parallel p(z \mid x)) \quad (12)$$

where KL divergence is defined as

$$\text{KL}(q \parallel p) = E_q[\log q(z)] - E_q[\log p(z \mid x)] \quad (13)$$

# More about KL Divergence

The Kullback–Leibler divergence measure the difference between two distributions

$$\text{KL}(q(x)||p(x)) = \sum_q q(x) \log \frac{q(x)}{p(x)} \quad (14)$$

$$= E_q[\log q(x)] - E_q[\log p(x)] \quad (15)$$

- ▶  $\text{KL}(q||p) = 0$ , if  $q = p$
- ▶  $\text{KL}(q||p) \geq 0$

Does not solve the problem

$$\text{KL}(q\|p) = E_q[\log q(\mathbf{z})] - E_q[\log p(\mathbf{z} \mid \mathbf{x})] \quad (16)$$

Does not solve the problem

$$\text{KL}(q\|p) = E_q[\log q(z)] - E_q[\log p(z \mid \mathbf{x})] \quad (16)$$

One more step we need

$$\text{KL}(q\|p) = E_q[\log q(z)] - E_q[\log p(z, \mathbf{x})] + \log p(\mathbf{x}) \quad (17)$$

Does not solve the problem

$$\text{KL}(q\|p) = E_q[\log q(\mathbf{z})] - E_q[\log p(\mathbf{z} \mid \mathbf{x})] \quad (16)$$

One more step we need

$$\text{KL}(q\|p) = E_q[\log q(\mathbf{z})] - E_q[\log p(\mathbf{z}, \mathbf{x})] + \log p(\mathbf{x}) \quad (17)$$

Evidence lower bound

$$\text{ELBo} = E_q[\log p(\mathbf{z}, \mathbf{x})] - E_q[\log q(\mathbf{z})] \quad (18)$$



# Justification

Put parameters back and consider the log-likelihood

$$\log p(x; \theta) = \log \sum_z p(x, z; \theta)$$

(19)

# Justification

Put parameters back and consider the log-likelihood

$$\begin{aligned}\log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\ &= \log \sum_z \frac{p(x, z; \theta) q(z; \psi)}{q(z; \psi)}\end{aligned}$$

(19)

# Justification

Put parameters back and consider the log-likelihood

$$\begin{aligned}\log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\ &= \log \sum_z \frac{p(x, z; \theta) q(z; \psi)}{q(z; \psi)} \\ &\geq \sum_z q(z; \psi) \log \frac{p(x, z; \theta)}{q(z; \psi)}\end{aligned}\tag{19}$$

# Justification

Put parameters back and consider the log-likelihood

$$\begin{aligned}\log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\ &= \log \sum_z \frac{p(x, z; \theta) q(z; \psi)}{q(z; \psi)} \\ &\geq \sum_z q(z; \psi) \log \frac{p(x, z; \theta)}{q(z; \psi)} \\ &= \sum_z q(z; \psi) \log p(x, z; \theta) - \sum_z q(z; \psi) \log q(z; \psi)\end{aligned}\tag{19}$$

# Justification

Put parameters back and consider the log-likelihood

$$\begin{aligned}\log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\ &= \log \sum_z \frac{p(x, z; \theta) q(z; \psi)}{q(z; \psi)} \\ &\geq \sum_z q(z; \psi) \log \frac{p(x, z; \theta)}{q(z; \psi)} \\ &= \sum_z q(z; \psi) \log p(x, z; \theta) - \sum_z q(z; \psi) \log q(z; \psi) \\ &= E_q[\log p(z, x; \theta)] - E_q[\log q(z; \psi)]\end{aligned}\tag{19}$$

# Justification

Put parameters back and consider the log-likelihood

$$\begin{aligned}\log p(x; \theta) &= \log \sum_z p(x, z; \theta) \\&= \log \sum_z \frac{p(x, z; \theta) q(z; \psi)}{q(z; \psi)} \\&\geq \sum_z q(z; \psi) \log \frac{p(x, z; \theta)}{q(z; \psi)} \\&= \sum_z q(z; \psi) \log p(x, z; \theta) - \sum_z q(z; \psi) \log q(z; \psi) \\&= E_q[\log p(z, x; \theta)] - E_q[\log q(z; \psi)] \\&= E_q[\log p(z, x; \theta)] + H(q)\end{aligned}\tag{19}$$

# Mean-field Approximation

One special case of variational inference is called **mean field approximation**, which specifies a family of variational distribution, in which different latent variables  $z_i$  are independent with each other.

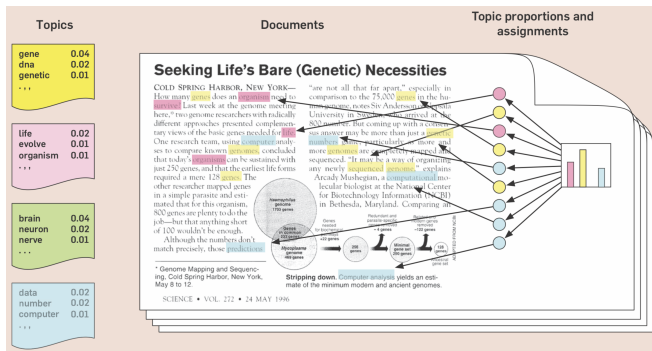
$$p(z; \psi) = \prod_i p(z_i; \psi_i) \quad (20)$$

## Example: Latent Dirichlet Allocation

---



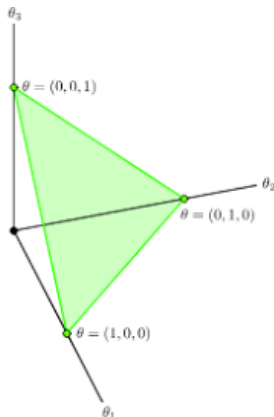
# Illustration



The basic idea is that a document is represented as a random *mixture* over latent topics, where each topic is characterized by a distribution over words. [Blei, 2012]

# Dirichlet Distribution

$$p(\boldsymbol{\theta}; \boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \theta_k^{\alpha_k - 1} \quad (21)$$



# Generative Story

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

where

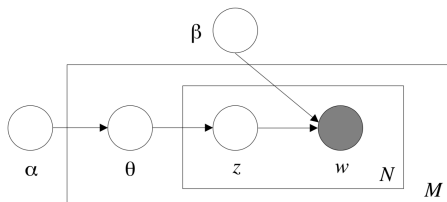
- ▶  $\theta \in \mathbb{R}^K$  is a  $K$ -dimensional random vector from the Dirichlet distribution with parameter  $\alpha$
- ▶  $\beta \in \mathbb{R}^{K \times V}$  is a matrix with  $\beta_{ij} = p(w_j = 1 \mid z_i = 1)$

# Joint Probability

For one document

$$p(\theta, z, d; \alpha, \beta) = p(\theta; \alpha) \prod_{n=1}^N \{p(z_n; \theta)p(w_n | z_n; \beta)\} \quad (22)$$

$M$  documents in a corpus



The key inference problem is to compute the posterior distribution of the hidden variable given a document

$$p(\boldsymbol{\theta}, z \mid d; \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{p(\boldsymbol{\theta}, z, d; \boldsymbol{\alpha}, \boldsymbol{\beta})}{p(d; \boldsymbol{\alpha}, \boldsymbol{\beta})} \quad (23)$$

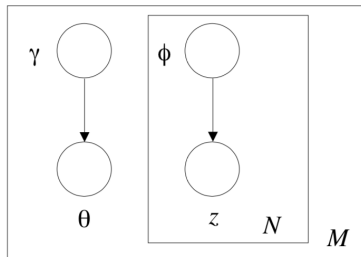
Recall that  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$  are the parameters of the original model.  $\boldsymbol{\theta}$  and  $z$  are latent variables.

# Variational Distribution

For one document

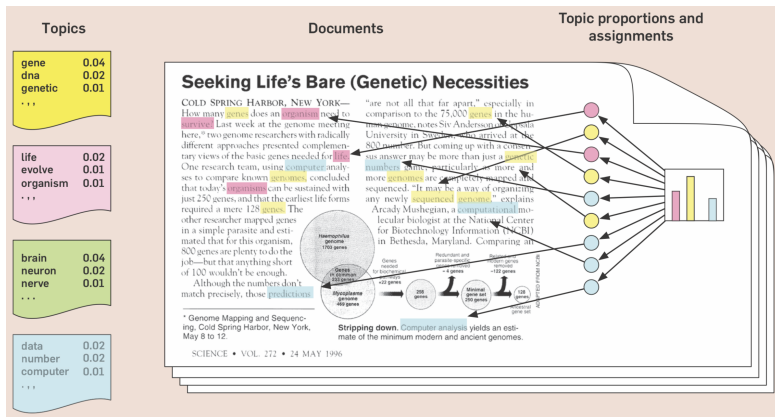
$$q(\theta, z; \gamma, \phi) = q(\theta; \gamma) \prod_n q(z_n; \phi) \quad (24)$$

$M$  documents in a corpus



$$\begin{aligned}\text{ELBo}_{\text{LDA}} = & E_q[\log p(\boldsymbol{\theta}; \boldsymbol{\alpha})] + E_q[\log p(\mathbf{z}; \boldsymbol{\theta}) \\ & + E_q[\log p(\mathbf{w} \mid \mathbf{z}; \boldsymbol{\beta})] \\ & - E_q[\log q(\boldsymbol{\theta}; \boldsymbol{\gamma})] - E_q[\log q(\mathbf{z}; \boldsymbol{\phi})]\end{aligned}\tag{25}$$

As shown in [Blei et al., 2003], every item in Eq. 25 has an analytic form, therefore we can have a closed form solution.





# Summary

1. Latent Variable Models
2. Variational Inference
3. Example: Latent Dirichlet Allocation

# Reference



Bishop, C. M. (2006).  
*Pattern Recognition and Machine Learning*.  
Springer-Verlag.



Blei, D. M. (2012).  
Probabilistic topic models.  
*Communications of the ACM*, 55(4):77–84.



Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003).  
Latent dirichlet allocation.  
*Journal of machine Learning research*, 3(Jan):993–1022.