

# CS 6501 Natural Language Processing

## Latent Variable Models

---

Yangfeng Ji

October 31, 2019

Department of Computer Science  
University of Virginia



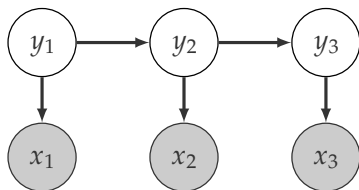
ENGINEERING

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

## Latent Variable Models

---

## Hidden Markov Models



# Gaussian Mixture Models

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

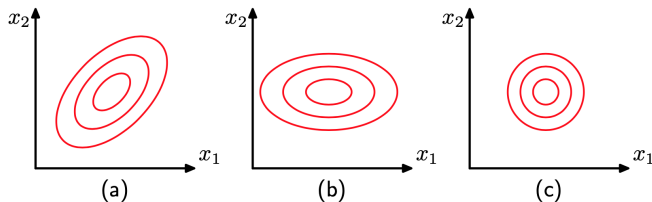
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad (1)$$

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 Distribution

$$\mathcal{N}(x \mid \mu, \Sigma) \propto \exp \left\{ -\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu) \right\} \quad (2)$$



► Mean:  $\mu$

► Covariance

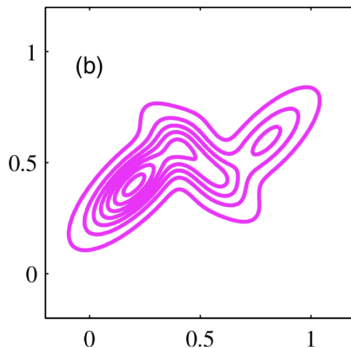
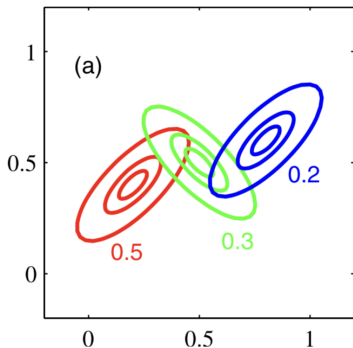
(a) : general form

(b) :  $\Sigma = \text{diag}(\sigma_i^2)$

(c) :  $\Sigma = \sigma^2 I$

# Gaussian Mixture Models: Example

$$p(\mathbf{x}) = \sum_{k=1}^3 \pi_k \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad (3)$$



[Bishop, 2006]

# GMM as a Latent Variable Model

Given a GMM

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad (4)$$

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

$$\mathbf{z} = [0, \dots, 0, 1, 0, \dots, 0] \quad (5)$$

- ▶ the probability of each component of  $\mathbf{z}$ ,  $z^{(k)}$ , is defined as

$$p(z^{(k)} = 1) = \pi_k \quad (6)$$



For each  $z_k$

$$p(z^{(k)} = 1) = \pi_k \quad (7)$$

Overall,

$$p(z) = \prod_{k=1}^K \pi_k^{z^{(k)}} \quad (8)$$

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

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

►  $p(\mathbf{x} \mid \mathbf{z}^{(k)} = 1)$

$$p(\mathbf{x} \mid \mathbf{z}^{(k)} = 1) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (9)$$

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

►  $p(\mathbf{x} \mid \mathbf{z}^{(k)} = 1)$

$$p(\mathbf{x} \mid \mathbf{z}^{(k)} = 1) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (9)$$

►  $p(\mathbf{x} \mid \mathbf{z})$

$$p(\mathbf{x} \mid \mathbf{z}) = \sum_{k=1}^K \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z^{(k)}} \quad (10)$$

Joint probability of the observed variable  $\mathbf{x}$  and latent variable  $\mathbf{z}$

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) \quad (11)$$

$$= \prod_{k=1}^K \pi_k^{z^{(k)}} \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z^{(k)}} \quad (12)$$

Joint probability of the observed variable  $x$  and latent variable  $z$

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

$$= \prod_{k=1}^K \pi_k^{z^{(k)}} \mathcal{N}(x | \mu_k, \Sigma_k)^{z^{(k)}} \quad (12)$$

Marginal probability of  $x$

$$p(x) = \sum_z p(x, z) \quad (13)$$

$$= \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k) \quad (14)$$

# Graphical Representation

With  $N$  data points from the GMM, each  $x_n$  has a latent variable  $z_n$  associated with it

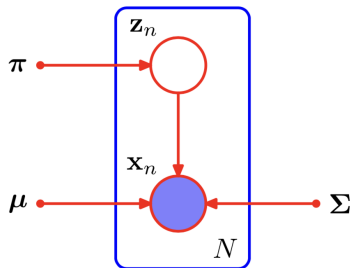
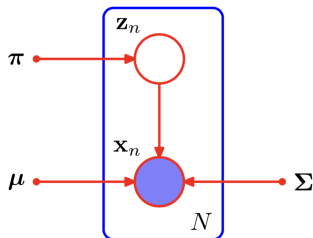


Figure: Graphical representation of GMM [Bishop, 2006]

# 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  $\{x_1, \dots, x_N\}$ , parameter estimation on a GMM is an iteration between the following two steps

1. Estimate  $p(z_n)$  for every  $x_n$
2. Estimate  $\{(\pi_k, \mu_k, \Sigma_k)\}_{k=1}^K$  based on  $\{p(z_n)\}$  and  $\{x_n\}$

Go back to step 1, until convergence



# Parameter Estimation

Given  $N$  data points  $\{x_1, \dots, x_N\}$ , parameter estimation on a GMM is an iteration between the following two steps

1. Estimate  $p(z_n)$  for every  $x_n$
2. Estimate  $\{(\pi_k, \mu_k, \Sigma_k)\}_{k=1}^K$  based on  $\{p(z_n)\}$  and  $\{x_n\}$

Go back to step 1, until convergence

## Comment

Similar to the  $K$ -means algorithm, with  $z_n$  as a random variable instead of a deterministic cluster assignment.

# Example

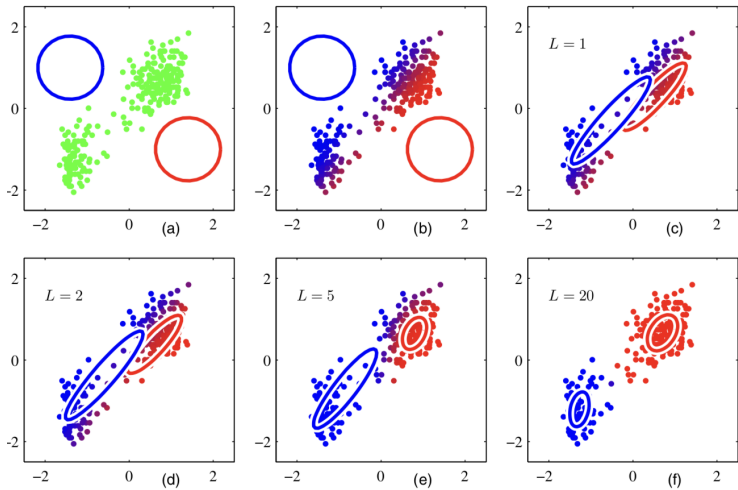
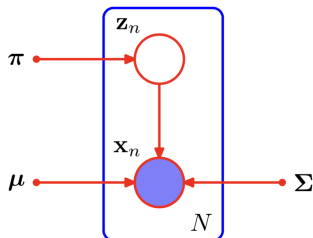


Figure: Illustration of the EM algorithm for GMM parameter estimation.



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

---

Recall the previous example

1. Estimate the probability of  $z$  using

$$p(z \mid x; \theta) = \frac{p(z, x)}{p(x)} \quad (15)$$

2. Estimate  $\theta$  using

$$\operatorname{argmax}_{\theta} \log p(z \mid x; \theta) \quad (16)$$

# Ideal Cases

Recall the previous example

1. Estimate the probability of  $z$  using

$$p(z \mid x; \theta) = \frac{p(z, x)}{p(x)} \quad (15)$$

2. Estimate  $\theta$  using

$$\operatorname{argmax}_{\theta} \log p(z \mid x; \theta) \quad (16)$$

## Key Requirement

$$p(x) = \sum_z p(x, z) \quad (17)$$

is tractable

However, the challenge comes from

$$p(\boldsymbol{x}) = \sum_{\boldsymbol{z}} p(\boldsymbol{x}, \boldsymbol{z}) \quad (18)$$

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

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

$$\tilde{q}(z) = \underset{q(z) \in \mathbb{Q}}{\operatorname{argmin}} \operatorname{KL}(q(z) \parallel p(z \mid x)) \quad (19)$$

where KL divergence is defined as

$$\operatorname{KL}(q \parallel p) = E_q[\log q(z)] - E_q[\log p(z \mid x)] \quad (20)$$



# More about KL Divergence

The Kullback–Leibler divergence measure the difference between two distributions

$$\text{KL}(q(x)\|p(x)) = \sum_q q(x) \log \frac{q(x)}{p(x)} \quad (21)$$

$$= E_q[\log q(x)] - E_q[\log p(x)] \quad (22)$$

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

$$\text{KL}(q\|p) = E_q[\log q(\mathbf{z})] - E_q[\log p(\mathbf{z} \mid \mathbf{x})] \quad (23)$$

$$\text{KL}(q\|p) = E_q[\log q(\mathbf{z})] - E_q[\log p(\mathbf{z} \mid \mathbf{x})] \quad (23)$$

One more step we need

$$\text{KL}(q\|p) = E_q[\log q(\mathbf{z})] - E_q[\log p(\mathbf{z}, \mathbf{x})] + \log p(\mathbf{x}) \quad (24)$$

$$\text{KL}(q\|p) = E_q[\log q(z)] - E_q[\log p(z \mid x)] \quad (23)$$

One more step we need

$$\text{KL}(q\|p) = E_q[\log q(z)] - E_q[\log p(z, x)] + \log p(x) \quad (24)$$

Evidence lower bound

$$\text{ELBo} = E_q[\log p(z, x)] - E_q[\log q(z)] \quad (25)$$

Consider maximizing the log-likelihood of the observed variable

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

Consider maximizing the log-likelihood of the observed variable

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

Consider maximizing the log-likelihood of the observed variable

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

Consider maximizing the log-likelihood of the observed variable

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



Consider maximizing the log-likelihood of the observed variable

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

Consider maximizing the log-likelihood of the observed variable

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

# Mean-field Approximation

A **special** case of variational inference is called **mean field approximation**, in which different latent variables  $z_i$  are independent with each other

$$q(\mathbf{z}; \boldsymbol{\psi}) = \prod_i q(z_i; \psi_i) \quad (26)$$

A **special** case of variational inference is called **mean field approximation**, in which different latent variables  $z_i$  are independent with each other

$$q(\mathbf{z}; \boldsymbol{\psi}) = \prod_i q(z_i; \psi_i) \quad (26)$$

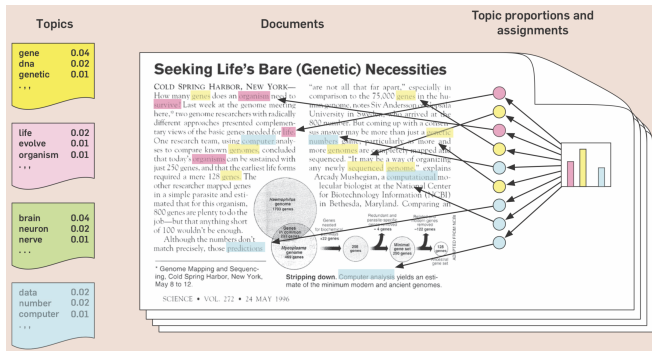
Requirements:

- ▶  $p(\mathbf{z}; \boldsymbol{\psi})$  is factorial
- ▶  $\psi_i$  can be computed with a close-form solution

## Example: Latent Dirichlet Allocation

---

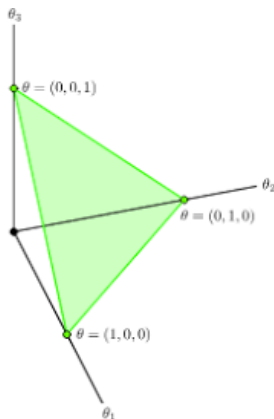
# Illustration



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

$$p(\boldsymbol{\theta}; \boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \theta_k^{\alpha_k - 1} \quad (27)$$



# 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  $\beta_{ij} = p(w_j = 1 \mid z_i = 1)$

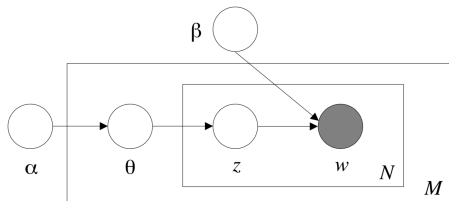


# Joint Probability

For one document

$$p(\theta, z, d; \alpha, \beta) = p(\theta; \alpha) \prod_{n=1}^N \{p(z_n; \theta)p(w_n | z_n; \beta)\} \quad (28)$$

$M$  documents in a corpus



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

$$p(\theta, z \mid d; \alpha, \beta) = \frac{p(\theta, z, d; \alpha, \beta)}{p(d; \alpha, \beta)} \quad (29)$$

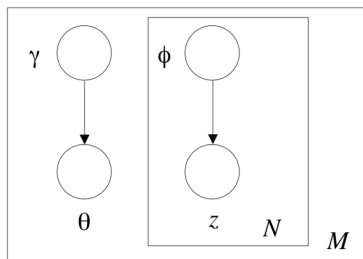
Recall that  $\alpha$  and  $\beta$  are the parameters of the original model.  $\theta$  and  $z$  are latent variables.

# Variational Distribution

For one document

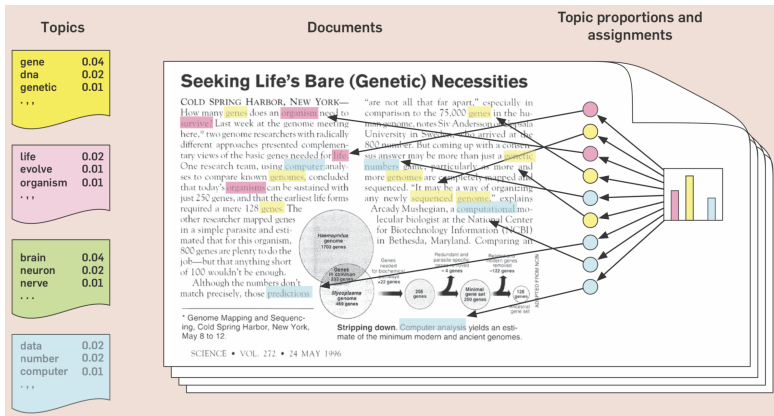
$$q(\theta, z; \gamma, \phi) = q(\theta; \gamma) \prod_n q(z_n; \phi) \quad (30)$$

$M$  documents in a corpus



$$\begin{aligned}\text{ELBo}_{\text{LDA}} = & E_q[\log p(\boldsymbol{\theta}; \boldsymbol{\alpha})] + E_q[\log p(\mathbf{z}; \boldsymbol{\theta}) \\ & + E_q[\log p(\mathbf{w} \mid \mathbf{z}; \boldsymbol{\beta})] \\ & - E_q[\log q(\boldsymbol{\theta}; \boldsymbol{\gamma})] - E_q[\log q(\mathbf{z}; \boldsymbol{\phi})]\end{aligned}\tag{31}$$

As shown in [Blei et al., 2003], every item in Eq. 31 has an analytic form, therefore we can have a closed form solution.



[Blei, 2012]

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.