# Learning Algorithms
# for Gaussian Mixture Models

**Department of Intelligent Science and Technology**
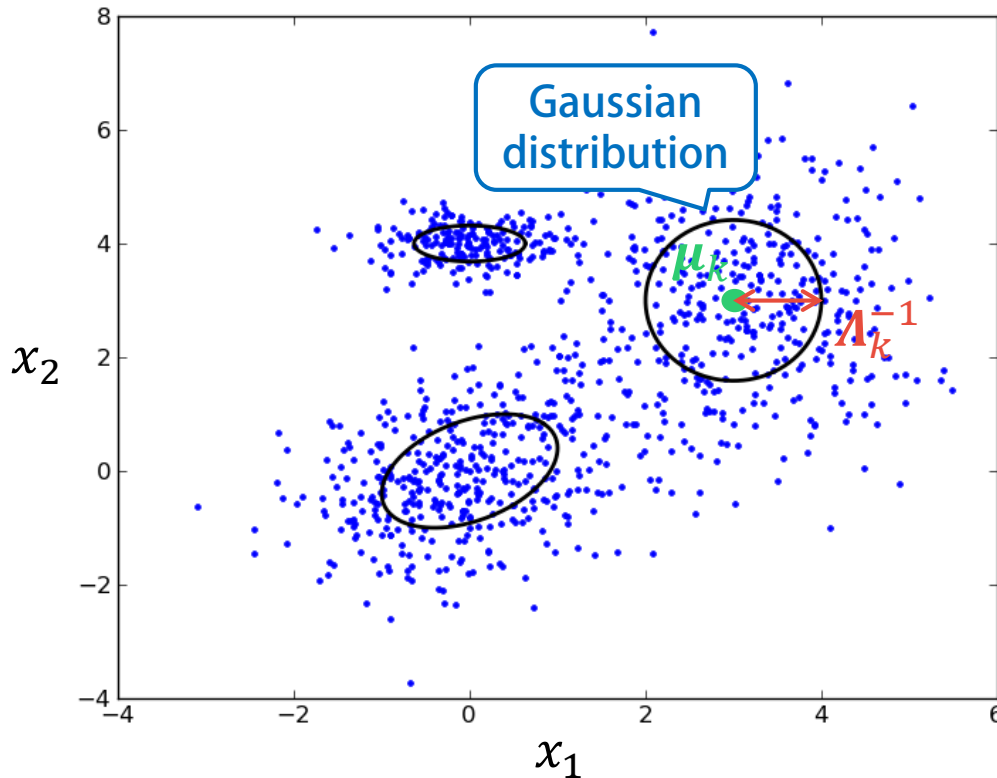**Graduate School of Informatics**
**Kyoto University**
**Kazuyoshi Yoshii**
yoshii@kuis.kyoto-u.ac.jp

# The Gaussian Mixture Model

- The GMM is used for representing how multi-dimensional vectors (e.g., feature vectors) are distributed stochastically



Probability distribution:

$$p(\boldsymbol{x}) = \sum_{k=1}^{K} \pi_k N(\boldsymbol{x}|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})$$

Parameters to be estimated:

Mixing ratios

$$\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]$$

Mean vectors

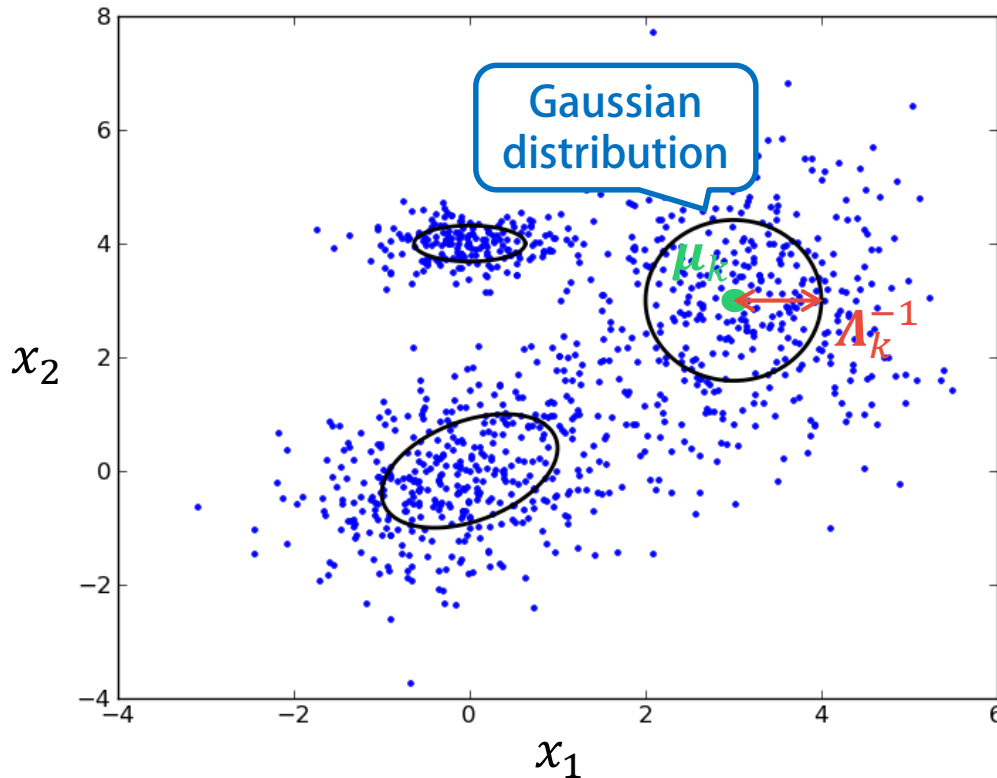$$\boldsymbol{\mu} = [\boldsymbol{\mu}_1, \cdots, \boldsymbol{\mu}_K]$$

Precision matrices

$$\boldsymbol{\Lambda} = [\boldsymbol{\Lambda}_1, \cdots, \boldsymbol{\Lambda}_K]$$

- ## The GMM is a probabilistic model for clustering
  - Each vector (sample) exclusively belongs to one of $K$ classes



Probability distribution:

$$p(\boldsymbol{x}) = \sum_{k=1}^{K} \pi_k N(\boldsymbol{x}|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})$$

Generative story:

Draw a latent variable

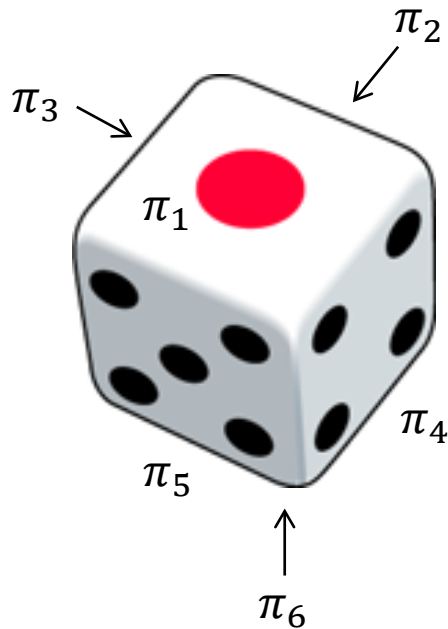$$\boldsymbol{z}_n \sim \text{Categorical}(\boldsymbol{z}_n|\boldsymbol{\pi})$$

Draw an observed variable

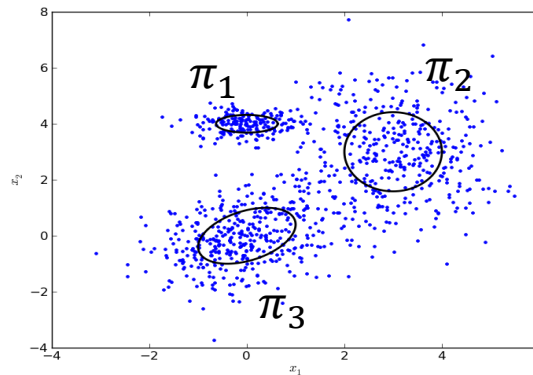$$\boldsymbol{x}_n \sim \prod_{k=1}^{K} N(\boldsymbol{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})^{z_{nk}}$$

- Latent variables are categorical distributed
  - Draw each latent variable: $z_n \sim \text{Categorical}(z_n|\pi)$ $(\pi = [\pi_1, \cdots, \pi_K])$
  - Use an one-of-$K$ representation: $z_n = [z_{n1}, z_{n2}, z_{n3}, \cdots, z_{nK}]$



$\pi_2$

$\pi_3$

$\pi_1$

$\pi_4$

$\pi_5$

$\pi_6$

Suppose we cast a $K$-sided die defined by $\pi$

If we get "3" for the $n^{th}$ trial, we say $z_n = [0, 0, 1, 0, 0, 0]$

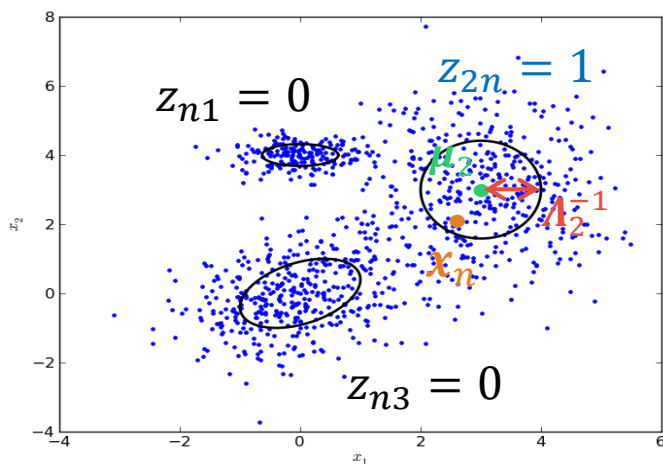Only one of the elements takes the value of 1



$\pi_1$

$\pi_2$

$\pi_3$

In the generative story of GMM, a class to which each sample belongs is stochastically determined by casting the die

- Observed variables are Gaussian distributed
    - Draw each observed variable: $x_n \sim \prod_{k=1}^{K} N(x_n | \mu_k, \Lambda_k^{-1})^{z_{nk}}$
    - Use the $k^{th}$ Gaussian distribution when $z_{nk} = 1$

Expand the product:

$$x_n \sim \prod_{k=1}^{3} N(x_n | \mu_k, \Lambda_k^{-1})^{z_{nk}} = N(x_n | \mu_1, \Lambda_1^{-1})^{z_{n1}} N(x_n | \mu_2, \Lambda_2^{-1})^{z_{n2}} N(x_n | \mu_3, \Lambda_3^{-1})^{z_{n3}}$$

Suppose $z_n = [0, 1, 0]$



$$x_n \sim N(x_n | \mu_2, \Lambda_2^{-1})$$

The one-of-$K$ representation can be used as a class indicator (selector)

This makes the derivation of learning algorithms easy (explained later)

- **There are several kinds of $K$-dimensional values**
  - **Random variables**
    - **Mixing ratios:** $\boldsymbol{\pi} = [\pi_1, \pi_2, \cdots, \pi_k, \cdots, \pi_K]$
    - **Latent variables:** $\boldsymbol{z}_n = [z_{n1}, z_{n2}, \cdots, z_{nk}, \cdots, z_{nK}]$
  - **Categorical probabilities**
    - **Posteriors:** $\boldsymbol{\gamma}_n = [\gamma_{n1}, \gamma_{n2}, \cdots, \gamma_{nk}, \cdots, \gamma_{nK}]$

**The values sum to unity**

$$\sum_{k=1}^{K} \pi_k = 1 \qquad \sum_{k=1}^{K} z_{nk} = 1 \qquad \sum_{k=1}^{K} \gamma_{nk} = 1$$

$0 < \pi_k < 1$

Only one of the values is 1
The other values are 0

$0 < \gamma_{nk} < 1$

- Generative story of the GMM
  - Draw each latent variable: $z_n \sim \text{Categorical}(z_n | \pi)$
  - Draw each observed variable: $x_n \sim \prod_{k=1}^{K} N(x_n | \mu_k, \Lambda_k^{-1})^{z_{nk}}$

- Two major approaches

| | Maximum likelihood (ML) estimation | Bayesian estimation |
|---|---|---|
| Probabilistic model | $p(X, Z; \mu, \Lambda)$ $= p(X|Z; \mu, \Lambda)p(Z; \pi)$ | $p(X, Z, \mu, \Lambda)$ $= p(X|Z, \mu, \Lambda)p(Z, \pi)p(\pi, \mu, \Lambda)$ |
| Latent variables $Z$ | Posterior calculation $p(Z|X; \pi, \mu, \Lambda)$ | Posterior calculation $p(Z, \pi, \mu, \Lambda|X)$ |
| Parameters $\pi, \mu, \Lambda$ | Point estimation $\pi^*, \mu^*, \Lambda^* = \arg\max p(X; \pi, \mu, \Lambda)$ | |

- ## Visualize dependency structures
  - Nodes: random variables (shaded = observable)
  - Edges: conditional dependencies



Likelihood model

Bayesian model

# Maximum Likelihood Estimation of Finite Gaussian Mixture Models

**Expectation-Maximization Algorithm**

**$K$-means Algorithm (Hard EM)**

- Suppose we have unlabeled height data
  - We want to estimate
    - the averages $\mu$ and precisions $\Lambda$ of the heights of male and female
    - the ratios $\pi$ of male and female



Count

Gender labels were lost!
Observed variables (heights): $X$
Latent variables (genders): $Z$

Assumption:
The height of each gender is
Gaussian distributed

Height

- Suppose we have unlabeled height data
  - We want to estimate
    - the averages $\mu$ and variances $\Lambda$ of the heights of male and female
    - the ratios $\pi$ of male and female

Count

If we <u>know</u> the gender of each sample, it is easy to calculate the above values
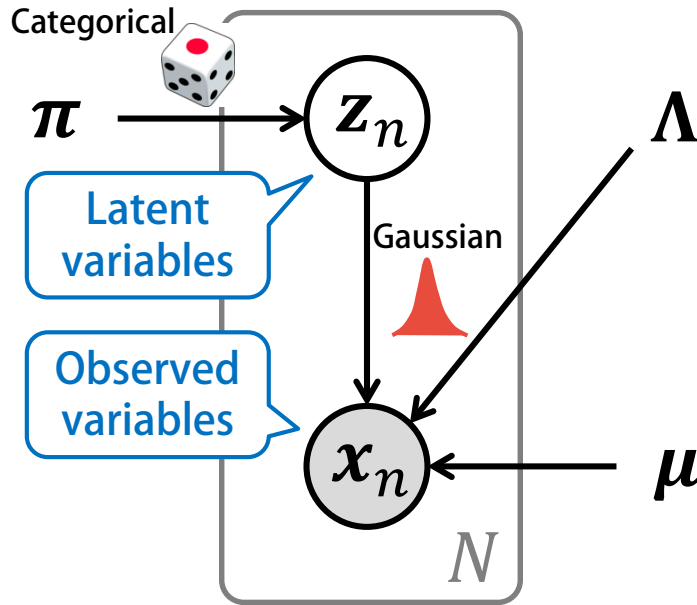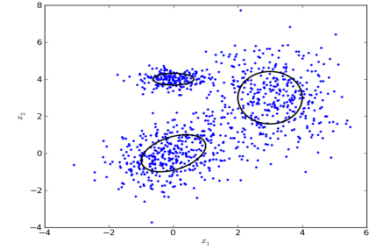
Our goal is unsupervised clustering of height data

Height

- Generative story of the GMM
  - Draw each latent variable: $z_n \sim \text{Categorical}(z_n | \pi)$
  - Draw each observed variable: $x_n \sim \prod_{k=1}^{K} N\left(x_n | \mu_k, \Lambda_k^{-1}\right)^{z_{nk}}$
- Two major approaches

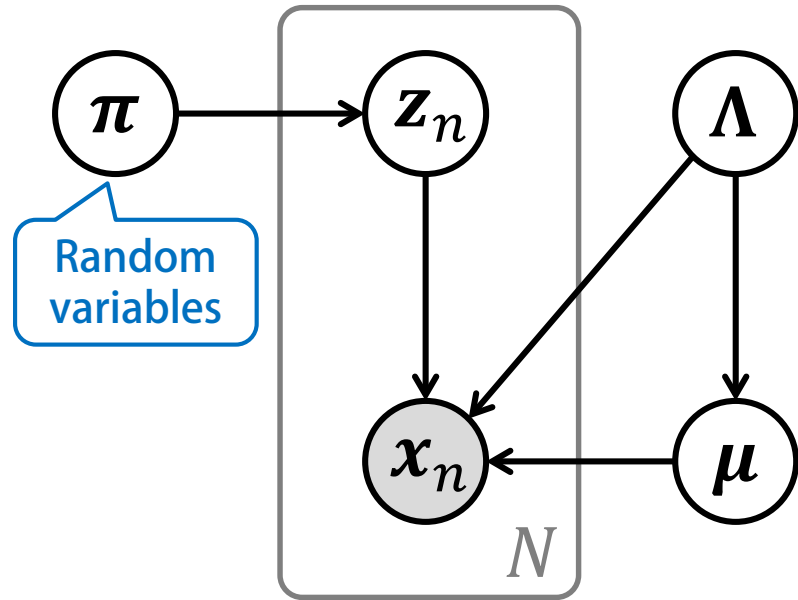| | Maximum likelihood (ML) estimation | Bayesian estimation |
|---|---|---|
| Probabilistic model | $p(X, Z; \mu, \Lambda)$ $= p(X|Z; \mu, \Lambda)p(Z; \pi)$ | $p(X, Z, \mu, \Lambda)$ $= p(X|Z, \mu, \Lambda)p(Z, \pi)p(\pi, \mu, \Lambda)$ |
| Latent variables $Z$ | Posterior calculation $p(Z|X; \pi, \mu, \Lambda)$ | Posterior calculation $p(Z, \pi, \mu, \Lambda | X)$ |
| Parameters $\pi, \mu, \Lambda$ | Point estimation $\pi^*, \mu^*, \Lambda^* = \text{argmax } p(X; \pi, \mu, \Lambda)$ | |

- Estimate the ratios, averages, and variances

$k = 1 \quad k = 2$

$\boldsymbol{x}_1 = 180cm \qquad \boldsymbol{z}_1 = [1, 0]$

$\boldsymbol{x}_2 = 170cm \qquad \boldsymbol{z}_2 = [0, 1]$

$\boldsymbol{x}_3 = 166cm \qquad \boldsymbol{z}_3 = [1, 0]$

$\boldsymbol{x}_4 = 175cm \qquad \boldsymbol{z}_4 = [1, 0]$

$\boldsymbol{x}_5 = 160cm \qquad \boldsymbol{z}_5 = [1, 0]$

$\boldsymbol{x}_6 = 155cm \qquad \boldsymbol{z}_6 = [0, 1]$

$\boldsymbol{x}_7 = 165cm \qquad \boldsymbol{z}_7 = [0, 1]$

$\boldsymbol{x}_8 = 162cm \qquad \boldsymbol{z}_8 = [1, 0]$

$\boldsymbol{x}_9 = 150cm \qquad \boldsymbol{z}_9 = [0, 1]$

**Sufficient statistics** for each class $k$ (male or female)

$$S_k[1] = \sum_{n=1}^{N} z_{nk} \quad \text{Count}$$

$$S_k[\boldsymbol{x}] = \sum_{n=1}^{N} z_{nk}\, \boldsymbol{x}_n \quad \text{Sum}$$

$$S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n=1}^{N} z_{nk}\, \boldsymbol{x}_n \boldsymbol{x}_n^T$$

Ratio: $\pi_k = \dfrac{S_k[1]}{S.[1]}$ 　　 Average: $\boldsymbol{\mu}_k = \dfrac{S_k[\boldsymbol{x}]}{S_k[1]}$ 　　 Variance: $\Lambda_k^{-1} = \dfrac{S_k[\boldsymbol{x}\boldsymbol{x}]}{S_k[1]} - \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T$

- Use posteriors instead of latent variables

$$k = 1 \quad k = 2$$

$$\boldsymbol{x}_1 = 180cm \qquad \boldsymbol{z}_1 = [?, ?]$$
$$\boldsymbol{x}_2 = 170cm \qquad \boldsymbol{z}_2 = [?, ?]$$
$$\boldsymbol{x}_3 = 166cm \qquad \boldsymbol{z}_3 = [?, ?]$$
$$\boldsymbol{x}_4 = 175cm \qquad \boldsymbol{z}_4 = [?, ?]$$
$$\boldsymbol{x}_5 = 160cm \qquad \boldsymbol{z}_5 = [?, ?]$$
$$\boldsymbol{x}_6 = 155cm \qquad \boldsymbol{z}_6 = [?, ?]$$
$$\boldsymbol{x}_7 = 165cm \qquad \boldsymbol{z}_7 = [?, ?]$$
$$\boldsymbol{x}_8 = 162cm \qquad \boldsymbol{z}_8 = [?, ?]$$
$$\boldsymbol{x}_9 = 150cm \qquad \boldsymbol{z}_9 = [?, ?]$$

$$k = 1 \quad k = 2$$

$$p(\boldsymbol{z}_1|\boldsymbol{X}) = [0.99, 0.01]$$
$$p(\boldsymbol{z}_2|\boldsymbol{X}) = [0.90, 0.10]$$
$$p(\boldsymbol{z}_3|\boldsymbol{X}) = [0.60, 0.40]$$
$$p(\boldsymbol{z}_4|\boldsymbol{X}) = [0.95, 0.05]$$
$$p(\boldsymbol{z}_5|\boldsymbol{X}) = [0.10, 0.90]$$
$$p(\boldsymbol{z}_6|\boldsymbol{X}) = [0.05, 0.95]$$
$$p(\boldsymbol{z}_7|\boldsymbol{X}) = [0.50, 0.50]$$
$$p(\boldsymbol{z}_8|\boldsymbol{X}) = [0.30, 0.70]$$
$$p(\boldsymbol{z}_9|\boldsymbol{X}) = [0.01, 0.99]$$

> Responsibility
> $\boldsymbol{\gamma}_n = [\gamma_{n1}, \gamma_{n2}]$

We cannot say $z_{nk} = 1$ for some $k$ with absolute certainty

To deal with uncertainty, we estimate the posterior of $z_{nk} = 1$

- Use posteriors instead of latent variables
  - Take into account the uncertainty of latent variables (genders)

<center>Genders known</center>

$$S_k[1] = \sum_{n=1}^{N} z_{nk} \qquad S_k[\boldsymbol{x}] = \sum_{n=1}^{N} z_{nk} \, \boldsymbol{x}_n$$

$$S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n=1}^{N} z_{nk} \, \boldsymbol{x}_n \boldsymbol{x}_n^T$$

<center>Genders unknown</center>

$$S_k[1] = \sum_{n=1}^{N} \gamma_{nk} \qquad S_k[\boldsymbol{x}] = \sum_{n=1}^{N} \gamma_{nk} \, \boldsymbol{x}_n$$

$$S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n=1}^{N} \gamma_{nk} \, \boldsymbol{x}_n \boldsymbol{x}_n^T$$

Ratio: $\pi_k^* = \dfrac{S_k[1]}{S_.[1]}$    Average: $\boldsymbol{\mu}_k^* = \dfrac{S_k[\boldsymbol{x}]}{S_k[1]}$    Variance: $\boldsymbol{\Lambda}_k^{-1\,*} = \dfrac{S_k[\boldsymbol{x}\boldsymbol{x}]}{S_k[1]} - \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T$

<center>How to estimate $z$ or $\gamma$</center>

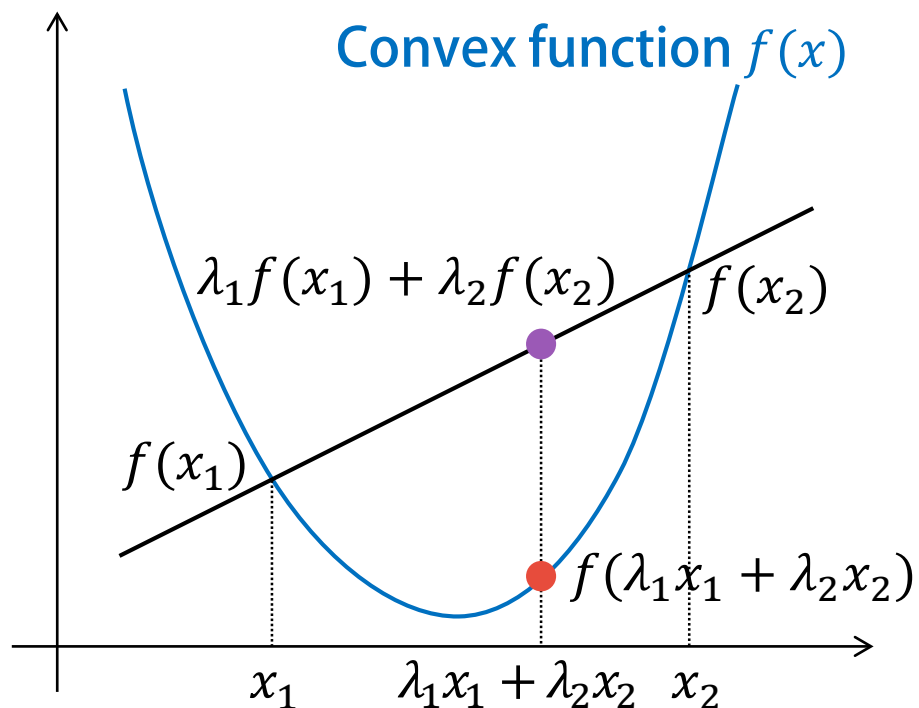| $K$-means algorithm (hard EM) (deterministic <u>hard</u> assignment) | EM algorithm (deterministic <u>soft</u> assignment) |

- Generative story of the GMM
  - Draw each latent variable: $z_n \sim \text{Categorical}(z_n|\pi)$
  - Draw each observed variable: $x_n \sim \prod_{k=1}^{K} N(x_n|\mu_k, \Lambda_k^{-1})^{z_{nk}}$

- Two major approaches

| | Maximum likelihood (ML) estimation | Bayesian estimation |
|---|---|---|
| Probabilistic model | $p(X, Z; \mu, \Lambda)$ $= p(X|Z; \mu, \Lambda)p(Z; \pi)$ | $p(X, Z, \mu, \Lambda)$ $= p(X|Z, \mu, \Lambda)p(Z, \pi)p(\pi, \mu, \Lambda)$ |
| Latent variables $Z$ | Posterior calculation $p(Z|X; \pi, \mu, \Lambda)$ | Posterior calculation $p(Z, \pi, \mu, \Lambda|X)$ |
| Parameters $\pi, \mu, \Lambda$ | Point estimation $\pi^*, \mu^*, \Lambda^* = \text{argmax}\, p(X; \pi, \mu, \Lambda)$ | |

- A basic inequality for convex functions
  - Forms the basis of the EM and VB algorithms

**Convex function** $f(x)$

$\lambda_1 f(x_1) + \lambda_2 f(x_2)$    $f(x_2)$

$f(x_1)$

$f(\lambda_1 x_1 + \lambda_2 x_2)$

$x_1$    $\lambda_1 x_1 + \lambda_2 x_2$    $x_2$

$$f\left(\sum_{k=1}^{K} \lambda_k x_k\right) \leq \sum_{k=1}^{K} \lambda_k f(x_k)$$

for auxiliary variables $\lambda$

such that $\sum_{k=1}^{K} \lambda_k = 1$

$$f\left(\int q(x)\, x\, dx\right) \leq \int q(x) f(x) dx$$

for auxiliary distribution $q(x)$

such that $\int q(x) dx = 1$

- Change the order of "sum" and "convex function"
  - Example: negative log of sum → sum of negative log

$$-\log\left(\sum_{k=1}^{K} x_k\right) = -\log\left(\sum_{k=1}^{K} \lambda_k \frac{x_k}{\lambda_k}\right) \leq -\sum_{k=1}^{K} \lambda_k \log\left(\frac{x_k}{\lambda_k}\right) \overset{\text{def}}{=} U(\boldsymbol{\lambda})$$

**Upper bound**

When does the equality holds true (when is $U(\boldsymbol{x}, \boldsymbol{\lambda})$ minimized)?

→ Optimization problem with a constraint

→ Method of Lagrange multipliers

$$\sum_{k=1}^{K} \lambda_k = 1$$

$$F(\boldsymbol{\lambda}) \overset{\text{def}}{=} U(\boldsymbol{\lambda}) + \omega\left(1 - \sum_{k=1}^{K} \lambda_k\right) \longrightarrow \frac{\partial F(\boldsymbol{\lambda})}{\partial \lambda_k} = -\log x_k + \log \lambda_k + 1 - \omega$$

**Equality condition**

Solving $\frac{\partial F(\boldsymbol{\lambda})}{\partial \lambda_k} = 0$, we get $\lambda_k = x_k e^{\omega-1}$ $\longrightarrow$ $e^{\omega-1} = \frac{1}{\sum_{k=1}^{K} x_k}$ $\longrightarrow$ $\lambda_k = \frac{x_k}{\sum_{k=1}^{K} x_k}$

- Change the order of "sum" and "convex function"
  - Example: negative log of sum → sum of negative log

$$-\log \int p(x,z)dz = {\color{red}-\log \int q(z)\frac{p(x,z)}{q(z)}dz} \leq {\color{purple}-\int q(z)\log\frac{p(x,z)}{q(z)}} \stackrel{\text{def}}{=} U(q(x))$$

Upper bound

When does the equality holds true (when is $U(q(x))$ minimized)?
→ Optimization problem with a constraint
→ Method of Lagrange multipliers

$$\sum_{k=1}^{K} q(x) = 1$$

$$F(q(x)) \stackrel{\text{def}}{=} U(q(x)) + \omega\left(1 - \int q(x)dx\right)$$ → Minimize as in the previous slide

Equality condition

$$q(z) = \frac{p(x,z)}{\int p(x,z)dz} = \frac{p(x,z)}{p(x)} = p(z|x)$$

- A deterministic algorithm for ML estimation
  - Suppose a probabilistic model $p(X, Z; \theta) = p(X|Z; \theta)p(Z; \theta)$
    - $X$: observed variables   $Z$: latent variables   $\theta$: parameters
  - We aim to get ML estimates $\theta^* = \operatorname{argmax} p(X; \theta)$  **Intractable!**

$$\log p(X; \theta) = \log \int p(X, Z; \theta) dZ$$

$$= \log \int q(Z) \frac{p(X, Z; \theta)}{q(Z)} dZ$$

Introduce an arbitrary distribution $q(Z)$ called a variational distribution

$$\geq \int q(Z) \log \frac{p(X, Z; \theta)}{q(Z)} dZ$$

**Jensen's inequality**

The equality holds true when $q^*(Z) = p(Z|X; \theta)$

$$= \int q(Z) \log p(X, Z; \theta) dZ - \int q(Z) \log q(Z)$$

**E-step**

**M-step** $\longrightarrow$ Maximize lower bound with respect to $\theta$

Hard EM: $q^*(Z) = \delta_{Z^*}(Z)$
$Z^* = \operatorname{argmax} p(Z|X; \theta)$

- **Iterate** E-step and M-step **alternately**
  - **E-step**: Calculate a posterior distribution over latent variables $Z$

$$x_1 = 180cm \qquad z_1 = [?, ?] \qquad\qquad \boldsymbol{\gamma}_1 = p(z_1|X) = [0.99, 0.01]$$
$$x_2 = 170cm \qquad z_2 = [?, ?] \qquad\qquad \boldsymbol{\gamma}_2 = p(z_2|X) = [0.90, 0.10]$$
$$x_3 = 166cm \qquad z_3 = [?, ?] \qquad\qquad \boldsymbol{\gamma}_3 = p(z_3|X) = [0.60, 0.40]$$

$$q^*(\mathbf{Z}) = p(\mathbf{Z}|X; \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{n=1}^{N} p(z_n|x_n; \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{n=1}^{N} \prod_{k=1}^{K} \gamma_{nk}^{z_{nk}}$$

**Responsibility**

$$q^*(z_{nk} = 1) = p(z_{nk} = 1|x_n; \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})$$

How well the sample $x_n$
is explained by each cluster
=
How likely the sample $x_n$
was to be generated
from each cluster

$$= \frac{p(x_n, z_{nk} = 1; \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})}{\sum_{k'=1}^{K} p(x_n, z_{nk'} = 1; \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})}$$

$$= \frac{\pi_k N(x_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})}{\sum_{k'=1}^{K} \pi_{k'} N(x_n|\boldsymbol{\mu}_{k'}, \boldsymbol{\Lambda}_{k'}^{-1})} = \gamma_{nk}$$

- **Iterate E-step and M-step alternately**
  - **M-step**: Update parameters $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}$
    - Calculate sufficient statistics

$$S_k[1] = \sum_{n=1}^{N} \gamma_{nk} \qquad S_k[\boldsymbol{x}] = \sum_{n=1}^{N} \gamma_{nk}\, \boldsymbol{x}_n$$

$$S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n=1}^{N} \gamma_{nk}\, \boldsymbol{x}_n \boldsymbol{x}_n^T$$

  - Estimate parameters

Ratio: $\pi_k^* = \dfrac{S_k[1]}{S.[1]}$　　Mean: $\boldsymbol{\mu}_k^* = \dfrac{S_k[\boldsymbol{x}]}{S_k[1]}$

Variance: $\boldsymbol{\Lambda}_k^{-1^*} = \dfrac{S_k[\boldsymbol{x}\boldsymbol{x}]}{S_k[1]} - \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T$

- **Iterate <u>M-step</u> and M-step alternately**
  - M-step: Update latent variables $Z$ — Hard assignment
    - $Z^* = \operatorname{argmax} p(Z|X; \pi, \mu, \Lambda) = \operatorname{argmax} p(X, Z; \pi, \mu, \Lambda)$
  - M-step: Update parameters $\pi, \mu, \Lambda$
    - $\pi^*, \mu^*, \Lambda^* = \operatorname{argmax} p(X|Z; \pi, \mu, \Lambda) = \operatorname{argmax} p(X, Z; \pi, \mu, \Lambda)$



If the all $\Lambda_k$'s are same, the hard EM for GMM reduces to the *k*-means algorithm

- A key difference lies in how to deal with uncertainty

| | $K$-means algorithm | EM algorithm |
|---|---|---|
| Latent variables $Z$ | Optimizing | Marginalizing out |
| Parameters $\pi, \mu, \Lambda$ | Optimizing | Optimizing |

$$S_k[1] = \sum_{n=1}^{N} z_{nk} \quad S_k[\boldsymbol{x}] = \sum_{n=1}^{N} z_{nk}\, \boldsymbol{x}_n \quad \Rightarrow \quad S_k[1] = \sum_{n=1}^{N} \gamma_{nk} \quad S_k[\boldsymbol{x}] = \sum_{n=1}^{N} \gamma_{nk}\, \boldsymbol{x}_n$$

$$S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n=1}^{N} z_{nk}\, \boldsymbol{x}_n\boldsymbol{x}_n^T \qquad S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n=1}^{N} \gamma_{nk}\, \boldsymbol{x}_n\boldsymbol{x}_n^T$$

Ratio: $\pi_k^* = \dfrac{S_k[1]}{S.[1]}$    Mean: $\boldsymbol{\mu}_k^* = \dfrac{S_k[\boldsymbol{x}]}{S_k[1]}$    Variance: $\Lambda_k^{-1^*} = \dfrac{S_k[\boldsymbol{x}\boldsymbol{x}]}{S_k[1]} - \boldsymbol{\mu}_k\boldsymbol{\mu}_k^T$

# Bayesian Estimation of
# Finite Gaussian Mixture Models

**(Collapsed) Gibbs Sampling**
**(Collapsed) Variational Bayes**

- Regard parameters as random variables
  - Introduce prior distributions on parameters
    - The Dirichlet distribution
      - A conjugate prior on categorical distributions
    - The Gaussian-Wishart distribution
      - A conjugate prior on Gaussian distributions



Gaussian-Wishart

- **Regard parameters as random variables**
  - Introduce prior distributions on parameters
  - Calculate posterior distributions on random variables

---

**Maximum likelihood estimation**

Latent variables: $p(Z|X; \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})$

Ratio: $\pi_k^* = \dfrac{S_k[1]}{S.[1]}$     Mean: $\boldsymbol{\mu}_k^* = \dfrac{S_k[x]}{S_k[1]}$     Variance: $\boldsymbol{\Lambda}_k^{-1^*} = \dfrac{S_k[xx]}{S_k[1]} - \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T$

---

**Bayesian estimation**

$\underset{\text{Likelihood}}{p(X|Z, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(Z|\boldsymbol{\pi})}$   x   $\underset{\text{Prior}}{p(\boldsymbol{\pi}) p(\boldsymbol{\mu}, \boldsymbol{\Lambda})}$   $\longrightarrow$   $\underset{\text{Posterior}}{p(Z, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}|X)}$

Bayes' theorem: $p(Z, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}|X) = \dfrac{p(X|Z, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(Z|\boldsymbol{\pi}) p(\boldsymbol{\pi}) p(\boldsymbol{\mu}, \boldsymbol{\Lambda})}{p(X)}$

- ## Widely used for mathematical convenience
  - The posterior $p(\boldsymbol{\theta}|X)$ takes the same form of the prior $p(\boldsymbol{\theta})$ for a particular type of the likelihood $p(X|\boldsymbol{\theta})$
    - $p(\boldsymbol{\pi}), p(\boldsymbol{\pi}|Z)$: Dirichlet $\qquad\qquad p(Z|\boldsymbol{\pi})$: Categorical
    - $p(\boldsymbol{\mu}, \boldsymbol{\Lambda}), p(\boldsymbol{\mu}, \boldsymbol{\Lambda}|X, Z)$: Gaussian-Wishart $\quad p(X|Z, \boldsymbol{\mu}, \boldsymbol{\Lambda})$: Gaussian

$\mathrm{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha})$ — Changing $\alpha$ from 0 to 2

2-dim. simplex
$\theta_1 + \theta_2 + \theta_3 = 1$

$\theta_3$

$\theta_2$

$\theta_1$

$\theta_2$

$\theta_3$

$\theta_1$

$\alpha = [6,2,2]$

$\alpha = [3,7,5]$

$\alpha = [2,3,4]$

$\alpha = [6,2,6]$

- Generative story of the GMM
  - Draw each latent variable: $z_n \sim \text{Categorical}(z_n|\pi)$
  - Draw each observed variable: $x_n \sim \prod_{k=1}^{K} N\left(x_n|\mu_k, \Lambda_k^{-1}\right)^{z_{nk}}$

- Two major approaches

| | Maximum likelihood (ML) estimation | Bayesian estimation |
|---|---|---|
| Probabilistic model | $p(X, Z; \mu, \Lambda)$ $= p(X|Z; \mu, \Lambda)p(Z; \pi)$ | $p(X, Z, \mu, \Lambda)$ $= p(X|Z, \mu, \Lambda)p(Z, \pi)p(\pi, \mu, \Lambda)$ |
| Latent variables $Z$ | Posterior calculation $p(Z|X; \pi, \mu, \Lambda)$ | Posterior calculation $p(Z, \pi, \mu, \Lambda|X)$ |
| Parameters $\pi, \mu, \Lambda$ | Point estimation $\pi^*, \mu^*, \Lambda^* = \text{argmax}\, p(X; \pi, \mu, \Lambda)$ | |

- Estimate the ratios, averages, and variances

$$k = 1 \quad k = 2$$

$x_1 = 180cm$  $z_1 = [1, 0]$

$x_2 = 170cm$  $z_2 = [0, 1]$

$x_3 = 166cm$  $z_3 = [1, 0]$

$x_4 = 175cm$  $z_4 = [1, 0]$

$x_5 = 160cm$  $z_5 = [1, 0]$

$x_6 = 155cm$  $z_6 = [0, 1]$

$x_7 = 165cm$  $z_7 = [0, 1]$

$x_8 = 162cm$  $z_8 = [1, 0]$

$x_9 = 150cm$  $z_9 = [0, 1]$

Sufficient statistics for each cluster $k$ (male or female)

$$S_k[1] = \sum_{n=1}^{N} z_{nk} \quad \text{Count}$$

$$S_k[\boldsymbol{x}] = \sum_{n=1}^{N} z_{nk}\, \boldsymbol{x}_n \quad \text{Sum}$$

$$S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n=1}^{N} z_{nk}\, \boldsymbol{x}_n \boldsymbol{x}_n^T$$

How to calculate the posterior distribution $p(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda} | \boldsymbol{X}, \boldsymbol{Z})$ ?

- ## Calculate a posterior distribution on parameters $\boldsymbol{\pi}$
  - ### The generative story
    - Prior: $\boldsymbol{\pi} \sim \mathrm{Dir}(\boldsymbol{\alpha}_0)$
    - Likelihood: $\boldsymbol{z}_n \sim \mathrm{Categorical}(\boldsymbol{z}_n | \boldsymbol{\pi})$

Posterior count | Prior count | Actual count

$$\alpha_k = \alpha_{0k} + S_k[1]$$

$$p(\boldsymbol{\pi}) = \mathrm{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha}_0) = \frac{\Gamma\left(\sum_{k=1}^{K} \alpha_{0k}\right)}{\prod_{k=1}^{K} \Gamma(\alpha_{0k})} \prod_{k=1}^{K} \pi_k^{\alpha_{0k}-1}$$

Bayes' theorem:

$$p(\boldsymbol{\pi}|\boldsymbol{Z})$$

$$p(\boldsymbol{Z}|\boldsymbol{\pi}) = \prod_{n=1}^{N} \mathrm{Categorical}(\boldsymbol{z}_n|\boldsymbol{\pi}) = \prod_{n=1}^{N}\prod_{k=1}^{K} \pi_k^{z_{nk}}$$

$$= \frac{p(\boldsymbol{Z}|\boldsymbol{\pi})p(\boldsymbol{\pi})}{p(\boldsymbol{Z})}$$

$$\propto p(\boldsymbol{Z}|\boldsymbol{\pi})p(\boldsymbol{\pi})$$

$$p(\boldsymbol{\pi}|\boldsymbol{Z}) = \mathrm{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha}) \propto \frac{\Gamma\left(\sum_{k=1}^{K} \alpha_{0k}\right)}{\prod_{k=1}^{K} \Gamma(\alpha_{0k})} \prod_{k=1}^{K} \pi_k^{\alpha_{0k}+S_k[1]-1}$$

We do not need to directly calculate the normalizing factor

- ## Calculate a posterior distribution on parameters $\boldsymbol{\mu}, \boldsymbol{\Lambda}$
  - ### The generative story
    - Prior: $\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k \sim N(\boldsymbol{\mu}_k | \boldsymbol{m}_0, (\beta_0 \boldsymbol{\Lambda}_k)^{-1}) W(\boldsymbol{\Lambda}_k | \mathbf{W}_0, \nu_0)$
    - Likelihood: $\boldsymbol{x}_n \sim \prod_{k=1}^{K} N(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})^{z_{nk}}$

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{k=1}^{K} N(\boldsymbol{\mu}_k | \boldsymbol{m}_0, (\beta_0 \boldsymbol{\Lambda}_k)^{-1}) W(\boldsymbol{\Lambda}_k | \mathbf{W}_0, \nu_0)$$

$$p(\boldsymbol{X} | \boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{n=1}^{N} \prod_{k=1}^{K} N(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})^{z_{nk}}$$

Bayes' theorem

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \boldsymbol{X}, \boldsymbol{Z}) = \prod_{k=1}^{K} N(\boldsymbol{\mu}_k | \boldsymbol{m}_k, (\beta_k \boldsymbol{\Lambda}_k)^{-1}) W(\boldsymbol{\Lambda}_k | \mathbf{W}_k, \nu_k)$$

Posterior count | Prior count | Actual count

$$\beta_k = \beta_0 + S_k[1]$$

$$\boldsymbol{m}_k = \frac{\beta_0 \boldsymbol{m}_k + S_k[\boldsymbol{x}]}{\beta_0 + S_k[1]}$$

$$\nu_k = \nu_0 + S_k[1]$$

$$\boldsymbol{W}_k^{-1} = \boldsymbol{W}_0^{-1} + \beta_0 \boldsymbol{m}_0 \boldsymbol{m}_0^T + S_k[\boldsymbol{x}\boldsymbol{x}^T] - \beta_k \boldsymbol{m}_k \boldsymbol{m}_k^T$$

- Use posteriors instead of latent variables

$$k = 1 \quad k = 2 \qquad\qquad k = 1 \quad k = 2$$

$x_1 = 180cm \qquad z_1 = [?, ?] \qquad\qquad p(z_1|X) = [0.99, 0.01]$

$x_2 = 170cm \qquad z_2 = [?, ?] \qquad\qquad p(z_2|X) = [0.90, 0.10]$

$x_3 = 166cm \qquad z_3 = [?, ?] \qquad\qquad p(z_3|X) = [0.60, 0.40]$

$x_4 = 175cm \qquad z_4 = [?, ?] \qquad\qquad p(z_4|X) = [0.95, 0.05]$

$x_5 = 160cm \qquad z_5 = [?, ?] \qquad\qquad p(z_5|X) = [0.10, 0.90]$

$x_6 = 155cm \qquad z_6 = [?, ?] \qquad\qquad p(z_6|X) = [0.05, 0.95]$

$x_7 = 165cm \qquad z_7 = [?, ?] \qquad\qquad p(z_7|X) = [0.50, 0.50]$

$x_8 = 162cm \qquad z_8 = [?, ?] \qquad\qquad p(z_8|X) = [0.30, 0.70]$

$x_9 = 150cm \qquad z_9 = \qquad\qquad p(z_9|X) = [0.01, 0.99]$

> Responsibility
> $\boldsymbol{\gamma}_n = [\gamma_{n1}, \gamma_{n2}]$

We cannot say $z_{nk} = 1$ for some $k$ with absolute certainty

To deal with uncertainty, we estimate the posterior of $z_{nk} = 1$

- Use posteriors instead of latent variables
  - Take into account the uncertainty of latent variables (genders)

Hard assignment

$$S_k[1] = \sum_{n=1}^{N} z_{nk} \qquad S_k[\boldsymbol{x}] = \sum_{n=1}^{N} z_{nk}\,\boldsymbol{x}_n$$

$$S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n=1}^{N} z_{nk}\,\boldsymbol{x}_n\boldsymbol{x}_n^T$$

Soft assignment

$$S_k[1] = \sum_{n=1}^{N} \gamma_{nk} \qquad S_k[\boldsymbol{x}] = \sum_{n=1}^{N} \gamma_{nk}\,\boldsymbol{x}_n$$

$$S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n=1}^{N} \gamma_{nk}\,\boldsymbol{x}_n\boldsymbol{x}_n^T$$

How to estimate $z$ or $\gamma$

Gibbs sampling
(stochastic algorithm)

Variational Bayes
(deterministic algorithm)

- **Generative story of the GMM**
  - Draw each latent variable: $z_n \sim \text{Categorical}(z_n|\pi)$
  - Draw each observed variable: $x_n \sim \prod_{k=1}^K N\left(x_n|\mu_k, \Lambda_k^{-1}\right)^{z_{nk}}$

- **Two major approaches**

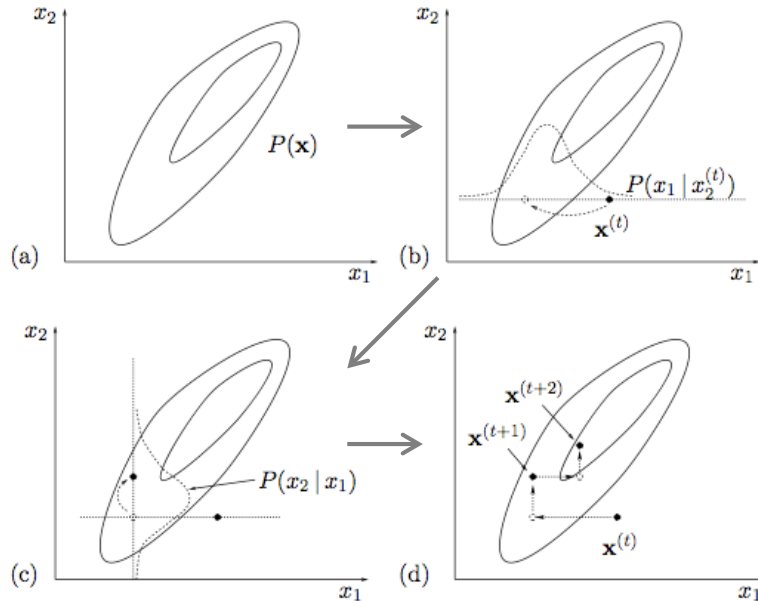|  | Maximum likelihood (ML) estimation | Bayesian estimation |
|---|---|---|
| Probabilistic model | $p(X, Z; \mu, \Lambda)$ $= p(X|Z; \mu, \Lambda)p(Z; \pi)$ | $p(X, Z, \mu, \Lambda)$ $= p(X|Z, \mu, \Lambda)p(Z, \pi)p(\pi, \mu, \Lambda)$ |
| Latent variables $Z$ | Posterior calculation $p(Z|X; \pi, \mu, \Lambda)$ | Posterior calculation $p(Z, \pi, \mu, \Lambda|X)$ |
| Parameters $\pi, \mu, \Lambda$ | Point estimation $\pi^*, \mu^*, \Lambda^* = \text{argmax}\, p(X; \pi, \mu, \Lambda)$ | |

- Choose an appropriate approach according to situations
  - Each approach has pros and cons
  - In general, Gibbs sampling is easy to implement

|  | Gibbs sampling | Variational Bayes |
|---|---|---|
| Convergence to true posterior | Yes | No |
| Judgment of convergence | Difficult | Easy |
| Convergence speed | Slow | Fast |
| Quality of estimation results | High | Moderate |

- A popular variant of Markov chain Monte Carlo (MCMC)
  - Generate random samples from a probability distribution $p(X) = \frac{f(X)}{Z}$ even if the normalizing factor $Z$ is intractable
  - The acceptance ratio is 100%



Objective: Generate independent samples from a probability distribution $p(X)$

1. Divide $X$ into several groups $X_1, \cdots, X_M$
2. for $t = 1 : T$
   for $m = 1 : M$
   Sample $X_m^{(t+1)}$

   $\sim p\left(X_m^{(t+1)} \middle| X_1^{(t+1)}, \cdots, X_{m-1}^{(t+1)}, X_{m+1}^{(t)}, \cdots, X_M^{(t)}\right)$
3. Pick up $X^{(t)}$ with a certain interval

This sampling needs to be done easily

- **Generate samples from** $p(\mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}|\mathbf{X})$

  

  - Divide $\{\mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}\}$ into $\{\mathbf{z}_1\}, \{\mathbf{z}_2\}, \cdots, \{\mathbf{z}_N\}, \{\boldsymbol{\pi}\}, \{\boldsymbol{\mu}, \boldsymbol{\Lambda}\}$
  - Iterate until convergence
    - for $n = 1:N$
      - Sample $\mathbf{z}_n \sim p(\mathbf{z}_n|\mathbf{X}, \mathbf{Z}_{-n}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = p(\mathbf{z}_n|\mathbf{x}_n, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})$
    - Sample $\boldsymbol{\pi} \sim p(\boldsymbol{\pi}|\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = p(\boldsymbol{\pi}|\mathbf{Z})$
    - Sample $\boldsymbol{\mu}, \boldsymbol{\Lambda} \sim p(\boldsymbol{\mu}, \boldsymbol{\Lambda}|\mathbf{X}, \mathbf{Z}, \boldsymbol{\pi}) = p(\boldsymbol{\mu}, \boldsymbol{\Lambda}|\mathbf{X}, \mathbf{Z})$

$$p(z_{nk} = 1|\mathbf{x}_n, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \frac{\pi_k N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)}{\sum_{k'=1}^{K} \pi_{k'} N(\mathbf{x}_n|\boldsymbol{\mu}_{k'}, \boldsymbol{\Lambda}_{k'})}$$

EM algorithm: soft assignment

Gibbs sampling: hard assignment

$$p(\boldsymbol{\pi}|\mathbf{Z}) = \mathrm{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha})$$

See "Posterior Calculation: Genders Known"

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda}|\mathbf{X}, \mathbf{Z}) = \prod_{k=1}^{K} N(\boldsymbol{\mu}_k|\mathbf{m}_k, (\beta_k \boldsymbol{\Lambda}_k)^{-1}) W(\boldsymbol{\Lambda}_k|\mathbf{W}_k, \nu_k)$$

- A Bayesian extension of the EM algorithm
  - We aim to approximate a true posterior $p(Z|X) = p(Z|X)/p(X)$ as a factorizable distribution $q(Z) = \prod_{m=1}^{M} q(Z_m)$

  **Intractable!**

$$\log p(X) = \log \int p(X, Z)dZ = \log \int q(Z)\frac{p(X, Z)}{q(Z)}dZ \geq \int q(Z)\log\frac{p(X, Z)}{q(Z)}dZ$$

$$= \int \left(\prod_{m=1}^{M} q(Z_m)\right)\left(\log p(X, Z) - \sum_{m=1}^{M} \log q(Z_m)\right)dZ_1 dZ_2 \cdots dZ_M$$

**Jensen's inequality**

$$= \sum_{m=1}^{M}\left(\int q(Z_m)\left(\int q(Z_{-m})\log p(X, Z)dZ_{-m}\right)dZ_m - \int q(Z_m)\log q(Z_m)dZ_m\right)$$

The lower bound is maximized when $\log q^*(Z_m) = \langle\log p(X, Z)\rangle_{q(Z_{-m})} + \text{const.}$

**The equality does NOT hold true!**

**VB-E step**   **VB-M step**

- VB just approximates a true posterior $p(\boldsymbol{Z}|\boldsymbol{X})$
  - The accuracy depends on how to factorize a variational posterior $q(\boldsymbol{Z})$

$$\log p(\boldsymbol{X}) = \int q(\boldsymbol{Z}) \log p(\boldsymbol{X}) dZ = \int q(\boldsymbol{Z}) \log \frac{q(\boldsymbol{Z})p(\boldsymbol{X}, \boldsymbol{Z})}{q(\boldsymbol{Z})p(\boldsymbol{Z}|\boldsymbol{X})} d\boldsymbol{Z}$$

$$= \int q(\boldsymbol{Z}) \log \frac{p(\boldsymbol{X}, \boldsymbol{Z})}{q(\boldsymbol{Z})} d\boldsymbol{Z} + \int q(\boldsymbol{Z}) \log \frac{q(\boldsymbol{Z})}{p(\boldsymbol{Z}|\boldsymbol{X})} d\boldsymbol{Z}$$

$$= \text{LowerBound}(q) + \text{KL}(q||p)$$

Kullback-Leibler (KL) divergence between
a variational posterior $q(\boldsymbol{Z})$
and a true posterior $p(\boldsymbol{Z}|\boldsymbol{X})$

Maximize = Minimize

The KD divergence is 0 when $q(\boldsymbol{Z}) = p(\boldsymbol{Z}|\boldsymbol{X})$ (intractable!)

If $q(\boldsymbol{Z})$ is assumed to be factorized, the KD divergence cannot be 0!

- Approximate a true posterior $p(\boldsymbol{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}|\boldsymbol{X})$
  - Assume a variational distribution $q(\boldsymbol{Z})q(\boldsymbol{\pi})q(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \approx p(\boldsymbol{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}|\boldsymbol{X})$
  - Iteratively update (optimize) each factor
    - VB-E step
      - $\log q^*(\boldsymbol{Z}) = \langle \log p(\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} + \text{const.}$

        $= \langle \log p(\boldsymbol{X}|\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{Z}|\boldsymbol{\pi}) \rangle_{q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} + \text{const.}$
    - VB-M step
      - $\log q^*(\boldsymbol{\pi}) = \langle \log p(\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} + \text{const.}$

        $= \langle \log p(\boldsymbol{Z}|\boldsymbol{\pi}) p(\boldsymbol{\pi}) \rangle_{q(\boldsymbol{Z})} + \text{const.}$
      - $\log q^*(\boldsymbol{\mu}, \boldsymbol{\Lambda}) = \langle \log p(\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{Z}, \boldsymbol{\pi})} + \text{const.}$

        $= \langle \log p(\boldsymbol{X}|\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{Z})} + \text{const.}$

Tractable posteriors: Use responsibilities instead of latent variables

- Formulate a full joint distribution

$$p(X, Z, \pi, \mu, \Lambda) = p(X|Z, \mu, \Lambda)p(Z|\pi)p(\pi)p(\mu, \Lambda)$$



$$p(X|Z, \mu, \Lambda) = \prod_{n=1}^{N}\prod_{k=1}^{K} N(x_n|\mu_k, \Lambda_k^{-1})^{z_{nk}}$$

$$p(Z|\pi) = \prod_{n=1}^{N} \text{Categorical}(z_n|\pi) = \prod_{n=1}^{N}\prod_{k=1}^{K} \pi_k^{z_{nk}}$$

Likelihood functions

$$p(\pi) = \text{Dir}(\pi|\alpha_0) = \frac{\Gamma\left(\sum_{k=1}^{K}\alpha_{0k}\right)}{\prod_{k=1}^{K}\Gamma(\alpha_{0k})}\prod_{k=1}^{K}\pi_k^{\alpha_{0k}-1}$$

$$p(\mu, \Lambda) = \prod_{k=1}^{K} N(\mu_k|m_0, (\beta_0\Lambda_k)^{-1})W(\Lambda_k|W_0, \nu_0)$$

Prior distributions

- **Invoke the updating formula of VB**
  - Take the expectation of the full joint probability distribution under "factorized" variational posteriors over other variables
  - Focus on only terms including $Z$
    (other terms can be absorbed into the normalization factor)

$$\log q^*(\boldsymbol{Z}) = \langle \log p(\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} + \text{const.}$$
$$= \langle \log p(\boldsymbol{X}|\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{Z}|\boldsymbol{\pi}) p(\boldsymbol{\pi}) p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} + \text{const.}$$
$$= \langle \log p(\boldsymbol{X}|\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{Z}|\boldsymbol{\pi}) \rangle_{q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} + \text{const.}$$

$$p(\boldsymbol{X}|\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{n=1}^{N} \prod_{k=1}^{K} N\left(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}\right)^{Z_{nk}}$$

$$p(\boldsymbol{Z}|\boldsymbol{\pi}) = \prod_{n=1}^{N} \text{Categorical}(\boldsymbol{z}_n | \boldsymbol{\pi}) = \prod_{n=1}^{N} \prod_{k=1}^{K} \pi_k^{Z_{nk}}$$

- Proceed the calculation according the updating rule

$$\langle \log p(\mathbf{Z}|\boldsymbol{\pi}) \rangle_{q(\boldsymbol{\pi})} = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \langle \log \pi_k \rangle_{q(\boldsymbol{\pi})}$$

$$\langle \log p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{\mu}, \boldsymbol{\Lambda})} = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \langle \log N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}) \rangle_{q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)}$$

$$\log q^*(\mathbf{Z}) = \langle \log p(\mathbf{Z}|\boldsymbol{\pi}) \rangle_{q(\boldsymbol{\pi})} + \langle \log p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{\mu}, \boldsymbol{\Lambda})} + \text{const.}$$

$$= \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \left( \langle \log \pi_k \rangle_{q(\boldsymbol{\pi})} + \langle \log N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}) \rangle_{q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)} \right) + \text{const.}$$

$$= \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \log \rho_{nk} + \text{const.}$$

- **Calculate the variational posterior over latent variables $Z$**
  - The normalization factor is automatically determined

$$\log q^*(\mathbf{Z}) = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \log \rho_{nk} + \text{const.}$$

> The distribution should be appropriately normalized

$$\gamma_{nk} = \frac{\rho_{nk}}{\sum_{k'=1}^{K} \rho_{nk'}}$$

$$\log q^*(\mathbf{Z}) = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \log \gamma_{nk}$$

$$q^*(\mathbf{Z}) = \prod_{n=1}^{N} \prod_{k=1}^{K} \gamma_{nk}^{z_{nk}} = \prod_{n=1}^{N} \text{Categorical}\left(\mathbf{z}_n \mid \boldsymbol{\gamma}_n\right)$$

> Latent variables are categorical distributed!

- ## Invoke the updating formula of VB

  - Take the expectation of the full joint probability distribution under "factorized" variational posteriors over other variables

  - Focus on only terms including $Z$
    (other terms can be absorbed into the normalization factor)

$$\log q^*(\boldsymbol{\pi}) = \langle \log p(\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} + \text{const.}$$
$$= \langle \log p(\boldsymbol{Z}|\boldsymbol{\pi}) p(\boldsymbol{\pi}) \rangle_{q(\boldsymbol{Z})} + \text{const.}$$
$$= \log p(\boldsymbol{\pi}) + \langle \log p(\boldsymbol{Z}|\boldsymbol{\pi}) \rangle_{q(\boldsymbol{Z})} + \text{const.}$$

$$\log q^*(\boldsymbol{\mu}, \boldsymbol{\Lambda}) = \langle \log p(\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{Z}, \boldsymbol{\pi})} + \text{const.}$$
$$= \langle \log p(\boldsymbol{X}|\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{Z})} + \text{const.}$$
$$= \log p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) + \langle \log p(\boldsymbol{X}|\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{Z})} + \text{const.}$$

> Bayesian estimation in simple conjugate models!
> (Use responsibilities $q(\boldsymbol{Z})$ instead of latent variables $\boldsymbol{Z}$)

- ## Calculate the variational posteriors over parameters $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}$
  - The posteriors take the same forms of the priors

$$S_k[1] = \sum_{n=1}^{N} \gamma_{nk} \quad S_k[\boldsymbol{x}] = \sum_{n=1}^{N} \gamma_{nk}\, \boldsymbol{x}_n \quad S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n=1}^{N} \gamma_{nk}\, \boldsymbol{x}_n \boldsymbol{x}_n^T$$

Sufficient statistics

$$p(\boldsymbol{\pi}) = \mathrm{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha}_0)$$
$$\downarrow$$
$$q^*(\boldsymbol{\pi}) = \mathrm{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha})$$

$$\alpha_k = \alpha_{0k} + S_k[1]$$

Posterior count · Prior count · Actual count

$$\beta_k = \beta_0 + S_k[1]$$

$$\boldsymbol{m}_k = \frac{\beta_0 \boldsymbol{m}_k + S_k[\boldsymbol{x}]}{\beta_0 + S_k[1]}$$

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{k=1}^{K} N(\boldsymbol{\mu}_k|\boldsymbol{m}_0, (\beta_0 \boldsymbol{\Lambda}_k)^{-1}) W(\boldsymbol{\Lambda}_k|\boldsymbol{W}_0, \nu_0)$$
$$\downarrow$$
$$q^*(\boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{k=1}^{K} N(\boldsymbol{\mu}_k|\boldsymbol{m}_k, (\beta_k \boldsymbol{\Lambda}_k)^{-1}) W(\boldsymbol{\Lambda}_k|\boldsymbol{W}_k, \nu_k)$$

$$\nu_k = \nu_0 + S_k[1]$$

$$\boldsymbol{W}_k^{-1} = \boldsymbol{W}_0^{-1} + \beta_0 \boldsymbol{m}_0 \boldsymbol{m}_0^T + S_k[\boldsymbol{x}\boldsymbol{x}^T] - \beta_k \boldsymbol{m}_k \boldsymbol{m}_k^T$$

- **Both methods have similar updating formulas**
  - EM: Using the values of parameters

$$\log p(\mathbf{Z}|\mathbf{X};\boldsymbol{\pi},\boldsymbol{\mu},\boldsymbol{\Lambda}) = \sum_{n=1}^{N}\sum_{k=1}^{K} z_{nk}\left(\log\pi_k + \log N(\mathbf{x}_n|\boldsymbol{\mu}_k,\boldsymbol{\Lambda}_k^{-1})\right) + \text{const.}$$

  - VB: Using the geometric means of parameters

$$\log q^*(\mathbf{Z}) = \sum_{n=1}^{N}\sum_{k=1}^{K} z_{nk}\left(\langle\log\pi_k\rangle_{q(\boldsymbol{\pi})} + \langle\log N(\mathbf{x}_n|\boldsymbol{\mu}_k,\boldsymbol{\Lambda}_k^{-1})\rangle_{q(\boldsymbol{\mu}_k,\boldsymbol{\Lambda}_k)}\right) + \text{const.}$$

$$\langle\log N(\mathbf{x}_n|\boldsymbol{\mu}_k,\boldsymbol{\Lambda}_k^{-1})\rangle_{q(\boldsymbol{\mu}_k,\boldsymbol{\Lambda}_k)} = -\frac{D}{2}\log(2\pi) + \frac{1}{2}\langle\log\boldsymbol{\Lambda}_k\rangle - \frac{1}{2}\left(\frac{D}{\beta_k^{-1}} + \nu_k(\mathbf{x}_m - \mathbf{m}_k)^T \mathbf{W}_k(\mathbf{x}_m - \mathbf{m}_k)\right)$$

$$\langle\log|\boldsymbol{\Lambda}_k|\rangle_{q(\boldsymbol{\pi})} = \sum_{d=1}^{D}\psi\left(\frac{c_k + 1 - d}{2}\right) + D\log 2 + \log|\mathbf{W}_k|$$

- ## The digamma function results in sparsifying effect

  - ▪ Example: Dirichlet distribution

    $$\boldsymbol{\pi} \sim \mathrm{Dir}(\boldsymbol{\alpha})$$

    - ◆ Mean

    $$\mathrm{E}[\pi_k] = \frac{\alpha_k}{\sum_{k'=1}^{K} \alpha_{k'}}$$
    $$= \exp\left(\log(\alpha_k) - \log\left(\sum_{k'=1}^{K} \alpha_{k'}\right)\right)$$

    - ◆ Geometric mean

    $$\mathrm{G}[\pi_k] = \exp(\mathrm{E}[\log \pi_k])$$
    $$= \exp\left(\psi(\alpha_k) - \psi\left(\sum_{k'=1}^{K} \alpha_{k'}\right)\right)$$

Small components tend to be degenerated in Bayesian mixture modeling

Discount!

$\exp\left(\log(\alpha_k)\right)$

$\exp\left(\psi(\alpha_k)\right)$

$\exp(\psi(0.1)) = 0.00003$
$\exp(\psi(0.5)) = 0.140$
$\exp(\psi(0.9)) = 0.470$

$\exp(\psi(1)) = 0.561$

$\exp(\psi(10)) = 9.504$
$\exp(\psi(100)) = 99.5004$
$\exp(\psi(1000)) = 999.500$

- **Both methods are based on similar updating formulas**
  - ▪ GS: <span style="color:#1a75c0">Stochastic hard</span> assignment

$$p(z_{nk} = 1 | \boldsymbol{x}_n, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \frac{\pi_k N(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)}{\sum_{k'=1}^{K} \pi_{k'} N(\boldsymbol{x}_n | \boldsymbol{\mu}_{k'}, \boldsymbol{\Lambda}_{k'})}$$

  - ▪ EM: <span style="color:#1a75c0">Deterministic soft</span> assignment

$$q^*(z_{nk} = 1 | \boldsymbol{x}_n, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \frac{\pi_k N(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)}{\sum_{k'=1}^{K} \pi_{k'} N(\boldsymbol{x}_n | \boldsymbol{\mu}_{k'}, \boldsymbol{\Lambda}_{k'})}$$

  - ▪ VB: <span style="color:#1a75c0">Deterministic soft</span> assignment

$$q^*(z_{nk} = 1) = \frac{G[\pi_k] G[N(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)]}{\sum_{k'=1}^{K} G[\pi_{k'}] G[N(\boldsymbol{x}_n | \boldsymbol{\mu}_{k'}, \boldsymbol{\Lambda}_{k'})]}$$

$$S_k[1] = \sum_{n=1}^{N} z_{nk}$$

$$S_k[\boldsymbol{x}] = \sum_{n=1}^{N} z_{nk} \boldsymbol{x}_n$$

$$S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n=1}^{N} z_{nk} \boldsymbol{x}_n \boldsymbol{x}_n^T$$

Replace $z_{nk}$ with $\gamma_{nk}$

- Reduce the number of variables for fast/better estimation
  - The parameters can be marginalized out due to conjugacy

$$p(X, Z, \pi, \mu, \Lambda) = p(X|Z, \mu, \Lambda)p(Z|\pi)p(\pi)p(\mu, \Lambda) \implies p(X|Z) = p(X|Z)p(Z)$$

$$p(Z|\pi) = \prod_{n=1}^{N} \text{Categorical}(z_n|\pi) = \prod_{n=1}^{N}\prod_{k=1}^{K} \pi_k^{z_{nk}}$$

Conjugacy holds true
(Dirichlet-Categorical)

$$p(\pi) = \text{Dir}(\pi|\alpha_0)$$

Marginalization over $\pi, \mu, \Lambda$ is analytically tractable!

$$p(X|Z, \mu, \Lambda) = \prod_{n=1}^{N}\prod_{k=1}^{K} N(x_n|\mu_k, \Lambda_k^{-1})^{z_{nk}}$$

Conjugacy holds true
(Gaussian-Wishart-Gaussian)

$$p(\mu, \Lambda) = \prod_{k=1}^{K} N(\mu_k|m_0, (\beta_0\Lambda_k)^{-1})W(\Lambda_k|W_0, \nu_0)$$

- Generate samples from $p(\boldsymbol{Z}|\boldsymbol{X})$
  - Divide $\{\boldsymbol{Z}\}$ into $\{\boldsymbol{z}_1\}, \{\boldsymbol{z}_2\}, \cdots, \{\boldsymbol{z}_N\}$
  - for $t = 1:T$
    - for $n = 1:N$
      - Sample $\boldsymbol{z}_n \sim p(\boldsymbol{z}_n|\boldsymbol{X}, \boldsymbol{Z}_{-n}) = p(\boldsymbol{z}_n|\boldsymbol{x}_n, \boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n})$

$$p(z_{nk} = 1|\boldsymbol{x}_n, \boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n}) \propto p(z_{nk} = 1, \boldsymbol{x}_n|\boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n})$$

$$= p(z_{nk} = 1|\boldsymbol{Z}_{-n})p(\boldsymbol{x}_n|z_{nk} = 1, \boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n})$$

$$= \int p(z_{nk} = 1|\boldsymbol{\pi})p(\boldsymbol{\pi}|\boldsymbol{Z}_{-n})d\boldsymbol{\pi} \int p(\boldsymbol{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)p(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k|\boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n})\, d\boldsymbol{\mu}_k d\boldsymbol{\Lambda}_k$$

$$= \frac{\alpha_k^{(-n)}}{\sum_{k'=1}^{K} \alpha_{k'}^{(-n)}} \mathrm{St}\left(\boldsymbol{x}_n\middle|\boldsymbol{m}_k^{(-n)}, \boldsymbol{L}_k^{(-n)}, \nu_k^{(-n)} + 1 - D\right)$$

Product of two predictive distributions

- Calculate predictive distributions
  - Marginalize likelihood functions under posteriors

$$\int \underbrace{p(z_{nk} = 1|\boldsymbol{\pi})}_{\text{Likelihood}} \underbrace{p(\boldsymbol{\pi}|\boldsymbol{Z}_{-n})}_{\text{Posterior}} d\boldsymbol{\pi} = \int \pi_k \text{Dir}\left(\boldsymbol{\pi}_k\middle|\boldsymbol{\alpha}^{(-n)}\right)d\boldsymbol{\pi} = \frac{\alpha_k^{(-n)}}{\sum_{k'=1}^K \alpha_{k'}^{(-n)}}$$

$$\int \underbrace{p(\boldsymbol{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)}_{\text{Likelihood}} \underbrace{p(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k|\boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n})}_{\text{Posterior}} d\boldsymbol{\mu}_k d\boldsymbol{\Lambda}_k$$

$$= \int N(\boldsymbol{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}) N\left(\boldsymbol{\mu}_k\middle|\boldsymbol{m}_k^{(-n)}, \left(\beta_k^{(-n)}\boldsymbol{\Lambda}_k\right)^{-1}\right) W\left(\boldsymbol{\Lambda}_k\middle|\boldsymbol{W}_k^{(-n)}, \nu_k^{(-n)}\right) d\boldsymbol{\mu}_k d\boldsymbol{\Lambda}_k$$

$$= \text{St}\left(\boldsymbol{x}_n|\boldsymbol{m}_k^{(-n)}, \boldsymbol{L}_k^{(-n)}, \nu_k^{(-n)} + 1 - D\right)$$

$$\boldsymbol{L}_k^{(-n)} = \frac{\nu_k^{(-n)} + 1 - D}{1 + \beta_k^{(-n)}} \boldsymbol{W}_k^{(-n)}$$

$$S_k[1] = \sum_{n' \neq N} z_{n'k} \quad S_k[\boldsymbol{x}] = \sum_{n' \neq n} z_{n'k}\boldsymbol{x}_{n'}$$

$$S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n' \neq n} z_{n'k}\boldsymbol{x}_{n'}\boldsymbol{x}_{n'}^T$$

- Approximate a posterior $p(\boldsymbol{Z}|\boldsymbol{X})$
  - Assume a variational distribution $\prod_{n=1}^{N} q(\boldsymbol{z}_n) \approx p(\boldsymbol{Z}|\boldsymbol{X})$
  - Iteratively update (optimize) each factor
    - CVB-E step: Invoke the updating formula of VB

$$\log q^*(\boldsymbol{z}_n) = \langle \log p(\boldsymbol{X}, \boldsymbol{Z}) \rangle_{q(\boldsymbol{Z}_{-n})} + \text{const.}$$

$$= \langle \log p(\boldsymbol{z}_n|\boldsymbol{X}, \boldsymbol{Z}_{-n}) p(\boldsymbol{X}|\boldsymbol{Z}_{-n}) p(\boldsymbol{Z}_{-n}) \rangle_{q(\boldsymbol{Z}_{-n})} + \text{const.}$$

$$= \langle \log p(\boldsymbol{z}_n|\boldsymbol{X}, \boldsymbol{Z}_{-n}) \rangle_{q(\boldsymbol{Z}_{-n})} + \text{const.}$$

$$p(z_{nk} = 1|\boldsymbol{X}, \boldsymbol{Z}_{-n}) \propto p(z_{nk} = 1, \boldsymbol{x}_n|\boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n})$$

$$= \frac{\alpha_k^{(-n)}}{\sum_{k'=1}^{K} \alpha_{k'}^{(-n)}} \text{St}\left(\boldsymbol{x}_n|\boldsymbol{m}_k^{(-n)}, \boldsymbol{L}_k^{(-n)}, v_k^{(-n)} + 1 - D\right)$$

Same as collapsed Gibbs sampling

- Calculate the variational posterior over latent variables $Z$
  - The normalization factor is automatically determined

$$\log q^*(z_{nk} = 1) = \langle \log p(\boldsymbol{z}_n | \boldsymbol{X}, \boldsymbol{Z}_{-n}) \rangle_{q(\boldsymbol{Z}_{-n})} + \text{const.}$$

$$= \left\langle \log \frac{\alpha_k^{(-n)}}{\sum_{k'=1}^{K} \alpha_{k'}^{(-n)}} + \log \text{St}\left(\boldsymbol{x}_n | \boldsymbol{m}_k^{(-n)}, \boldsymbol{L}_k^{(-n)}, v_k^{(-n)} + 1 - D\right) \right\rangle + \text{const.}$$

$$\approx \log \left\langle \alpha_k^{(-n)} \right\rangle - \log \sum_{k'=1}^{K} \left\langle \alpha_{k'}^{(-n)} \right\rangle$$

> **0-th order approximation (CVB0)**
> $$\mathrm{E}[\log x] \approx \log \mathrm{E}[x]$$

$$+ \log St\left(\boldsymbol{x}_n \middle| \left\langle \boldsymbol{m}_k^{(-n)} \right\rangle, \left\langle \boldsymbol{L}_k^{(-n)} \right\rangle, \left\langle v_k^{(-n)} \right\rangle + 1 - D\right) + \text{const.}$$

$$S_k[1] = \sum_{n' \neq N} \gamma_{n'k} \quad S_k[\boldsymbol{x}] = \sum_{n' \neq n} \gamma_{n'k} \boldsymbol{x}_{n'} \quad S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n' \neq n} \gamma_{n'k} \boldsymbol{x}_{n'} \boldsymbol{x}_{n'}^T$$

- Both methods are based on similar updating formulas
  - CGS: Stochastic hard assignment

$$p(z_{nk} = 1 | \boldsymbol{x}_n, \boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n}) = \frac{\alpha_k^{(-n)}}{\Sigma_{k'=1}^{K} \alpha_{k'}^{(-n)}} \, \text{St}\left(\boldsymbol{x}_n \middle| \boldsymbol{m}_k^{(-n)}, \boldsymbol{L}_k^{(-n)}, v_k^{(-n)} + 1 - D\right)$$

$$S_k[1] = \sum_{n' \neq N} z_{n'k} \quad S_k[\boldsymbol{x}] = \sum_{n' \neq n} z_{n'k} \boldsymbol{x}_{n'} \quad S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n' \neq n} z_{n'k} \boldsymbol{x}_{n'} \boldsymbol{x}_{n'}^T$$

  - CVB: Deterministic soft assignment

$$q(z_{nk} = 1) = \frac{\left\langle \alpha_k^{(-n)} \right\rangle}{\Sigma_{k'=1}^{K} \left\langle \alpha_{k'}^{(-n)} \right\rangle} \, \text{St}\left(\boldsymbol{x}_n \middle| \left\langle \boldsymbol{m}_k^{(-n)} \right\rangle, \left\langle \boldsymbol{L}_k^{(-n)} \right\rangle, \left\langle v_k^{(-n)} \right\rangle + 1 - D\right)$$

$$S_k[1] = \sum_{n' \neq N} \gamma_{n'k} \quad S_k[\boldsymbol{x}] = \sum_{n' \neq n} \gamma_{n'k} \boldsymbol{x}_{n'} \quad S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n' \neq n} \gamma_{n'k} \boldsymbol{x}_{n'} \boldsymbol{x}_{n'}^T$$

- All methods are based on similar updating formulas
  - GS: Stochastic hard assignment
    - $p(z_{nk} = 1|\boldsymbol{x}_n, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \propto \pi_k N(\boldsymbol{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$
  - CGS: Stochastic hard assignment
    - $p(z_{nk} = 1|\boldsymbol{x}_n, \boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n}) = \frac{\alpha_k^{(-n)}}{\sum_{k'=1}^K \alpha_{k'}^{(-n)}} \text{St}\left(\boldsymbol{x}_n \middle| \boldsymbol{m}_k^{(-n)}, \boldsymbol{L}_k^{(-n)}, \nu_k^{(-n)} + 1 - D\right)$
  - EM: Deterministic soft assignment
    - $q^*(z_{nk} = 1|\boldsymbol{x}_n, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \propto \pi_k N(\boldsymbol{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$
  - VB: Deterministic soft assignment
    - $q^*(z_{nk} = 1) \propto \text{G}[\pi_k]\text{G}[N(\boldsymbol{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)]$
  - CVB: Deterministic soft assignment
    - $q(z_{nk} = 1) = \frac{\left\langle \alpha_k^{(-n)} \right\rangle}{\sum_{k'=1}^K \left\langle \alpha_{k'}^{(-n)} \right\rangle} \text{St}\left(\boldsymbol{x}_n \middle| \left\langle \boldsymbol{m}_k^{(-n)} \right\rangle, \left\langle \boldsymbol{L}_k^{(-n)} \right\rangle, \left\langle \nu_k^{(-n)} \right\rangle + 1 - D\right)$

> All formulas are like:
> Mixing ratio
> ×
> Component
> distribution

- Learning algorithm can be categorized
  with respect to how to deal with uncertainty

| | | Latent variables $Z$ | | |
| --- | --- | --- | --- | --- |
| | | Point estimates | Posteriors | Sampled values |
| Parameters $\pi, \mu, \Lambda$ | Point estimates | $K$-means+ (maximization-maximization) | EM (expectation-maximization) | |
| | Posteriors | Bayesian $K$-means (maximization-expectation) | VB (expectation-expectation) | |
| | Sampled values | | | Gibbs sampling (sampling-sampling) |

- Implement basic functions for updating posteriors
  - Input: prior + statistics    Output: posterior

**posterior.h**

```cpp
void update_dirichlet
(mcl::Dirichlet& dirichlet,
 mcl::Dirichlet& dirichlet0,
 const std::vector<double>& s);
```

```cpp
void update_gaussian_wishart
(mcl::Gaussian& gaussian,
 mcl::Wishart& wishart,
 const mcl::Gaussian& gaussian0,
 const mcl::Wishart& wishart0,
 double s,
 const std::vector<double>& sx,
 const std::vector<double>& sxx);
```

```cpp
void update_student
(mcl::Student& student,
 const mcl::Gaussian& gaussian,
 const mcl::Wishart& wishart);
```

Predictive distribution
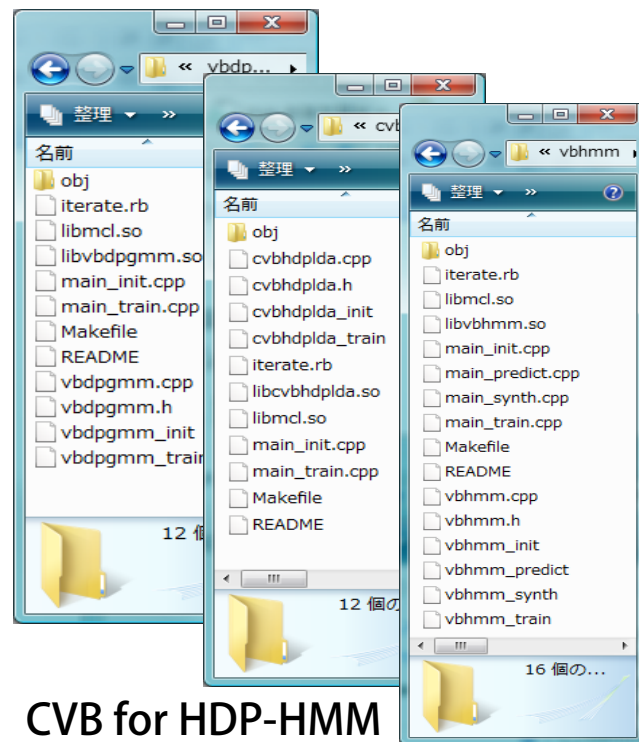(used for collapsed inference)

- Combine appropriate functions for your model
  - Use conjugate priors as much as possible

VB for DP-GMM



Library

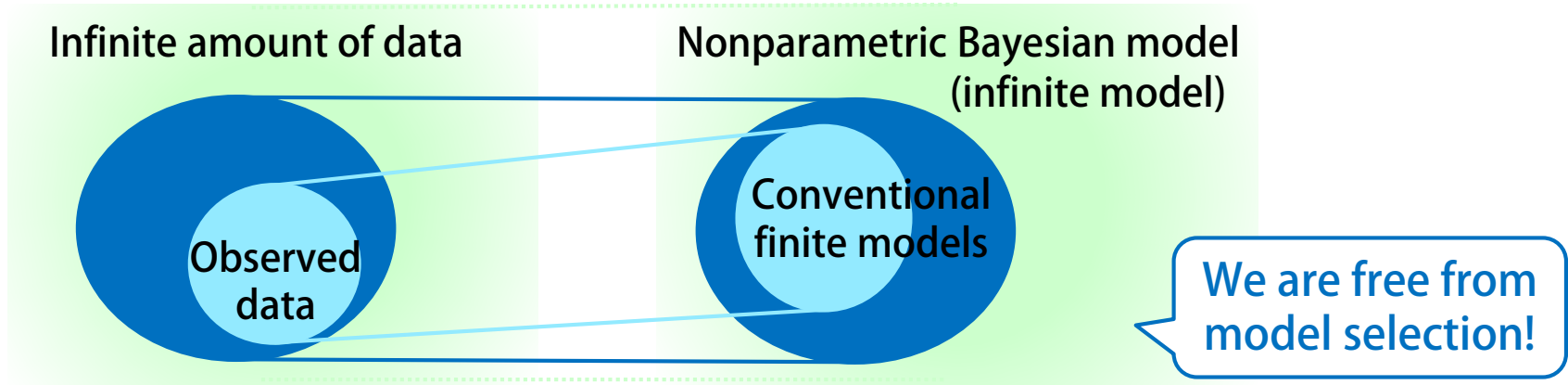CVB for HDP-HMM

VB for HMM

**MapReduce-type parallelization is easy**

- **Implement HTK-like commands**
  - `vbgmm_init [model.xml] [K]`
    - Make an initial model with K components
  - `vbgmm_train [model.xml] [data.csv] ([#iterations])`
    - Update the model using the data
    - Overwrite the model file
- **Parallelization based on boost::mpi**
  - MapReducing EM algorithm for Master-Slave architecture
    - E-step: *Master* distributes the data to *Slaves*
      - Each *Slave* calculates the responsibilities for the given data
    - M-step: *Master* gathers the responsibilities from *Slaves*
      - *Master* updates the posteriors

# Bayesian Estimation of
# Infinite Gaussian Mixture Models

## Collapsed Gibbs Sampling
## Variational Bayes

- Bayesian models with infinite complexity
  - "Nonparametric" means having an infinite number of parameters
  - Excellent generalization capability
    - If we have an infinite amount of data, all an infinite number of parameters are required
    - If we have a finite amount of data, only a finite subset is required

Infinite amount of data

Nonparametric Bayesian model
(infinite model)

Conventional finite models

Observed data

We are free from model selection!

- An infinite-dimensional prior distribution
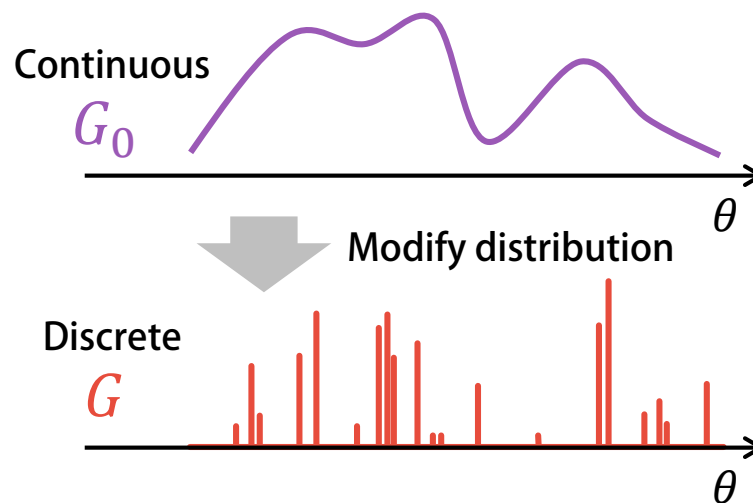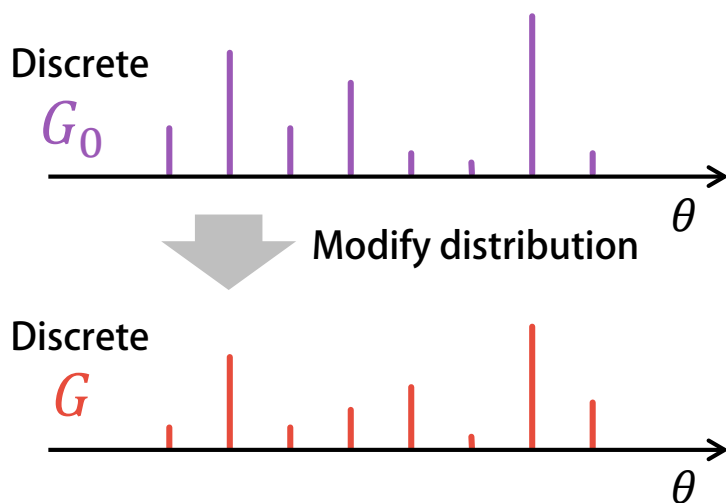  - Capable of generating infinite-dimensional distributions

$$G \sim \mathrm{DP}(\alpha, G_0)$$

Concentration parameter

Base measure

The DP can be explicitly rewritten as

$$G(\theta) = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}(\theta) \qquad \theta_k \sim G_0$$

Discrete $G_0$

$\theta$

Modify distribution

Discrete $G$

$\theta$

Continuous $G_0$

$\theta$

Modify distribution

Discrete $G$

$\theta$

- ## The DP always generates discrete distributions
  - ### The positions of "atoms" are shared with the discrete base measure

$$G(\theta) = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}(\theta) \quad \theta_k \sim G_0$$

Each $\theta_k$ is one of $\{\theta_1, \theta_2, \cdots, \theta_I\}$



Discrete $G_0$

Modify distribution

Discrete $G$

$\theta_1 \ \theta_2 \qquad \theta_i \qquad \theta_I \qquad \theta$

$$G(\theta) = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}(\theta)$$

$$= \sum_{i=1}^{I} \left( \sum_{k:\theta_k=\theta_i} \pi_k \right) \delta_{\theta_i}(\theta)$$

$I$-dimensional discrete distribution

- ## The DP always generates discrete distributions
  - The number of "atoms" are countably infinite

$$G(\theta) = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}(\theta) \quad \theta_k \sim G_0$$

$\theta_k$'s are almost surely disjoint
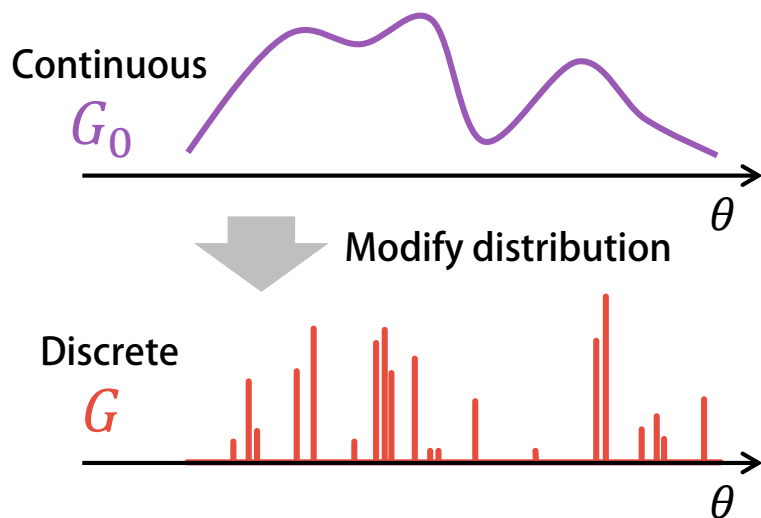
Continuous
$G_0$

$\theta$

Modify distribution

Discrete
$G$

$\theta$

If we use a continuous prior distribution as a base measure $G_0$, we can generate an infinite-dim. discrete distribution!

If $G_0$ is a Gaussian-Wishart distribution (the probability space is over $\theta = \{\boldsymbol{\mu}, \boldsymbol{\Lambda}\}$)

↓

$G$ consists of infinitely many Gaussians $\{\theta_1, \cdots, \theta_\infty\}$ with weights $\{\pi_1, \cdots, \pi_\infty\}$

- **Stochastically generate the weights** $\{\pi_1, \cdots, \pi_\infty\}$
  - a.k.a. Griffiths-Engen-McCloskey distribution

$$\boldsymbol{\pi} \sim \text{SBP}(\alpha) \text{ or } \text{GEM}(\alpha)$$

$$\downarrow$$

$$v_k \sim \text{Beta}(1, \alpha) \quad \pi_k = v_k \prod_{k'=1}^{k-1} (1 - v_{k'})$$

$\text{Beta}(\alpha, \beta)$



Generalization:

Pitman-Yor process

$v_k \sim \text{Beta}(1 - d, \alpha + dk)$

Beta two-parameter process

$\text{Beta}(\alpha, \beta)$

$\pi_1 \quad \pi_2 \quad \pi_3 \quad \cdots$

$v_1$

$1 - v_1$

$v_2$

$1 - v_2$

$v_3 \quad 1 - v_3$

- The concentration parameter controls the sparseness
  - The value of $\alpha$ is unknown $\rightarrow$ Introduce a hyper prior on $\alpha$

Continuous $G_0$

$v_k \sim \mathrm{Beta}(1, \alpha)$

$\pi_1$  $\pi_2$  $\pi_3$  ...

Large $\alpha$

Small $\alpha$

$\theta$

$v_k \sim \mathrm{Beta}(1, \alpha)$

$\pi_1$  $\pi_2$  $\pi_3$ ...

Discrete $G$

$\theta$

Discrete $G$

$\theta$

Many $\theta$'s are likely to be used for representing the data

Fewer $\theta$'s are likely to be used for representing the data

Assume $\alpha \sim \mathrm{Gamma}(a, b)$ for taking into account uncertainty

- ## Generate infinitely many Gaussians using a DP

$$G(\boldsymbol{\mu}, \boldsymbol{\Lambda}) = \sum_{k=1}^{\infty} \pi_k \delta_{\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k}(\boldsymbol{\mu}, \boldsymbol{\Lambda})$$

$$\boldsymbol{\pi} \sim \mathrm{SBP}(\alpha) \quad \text{SBP prior}$$

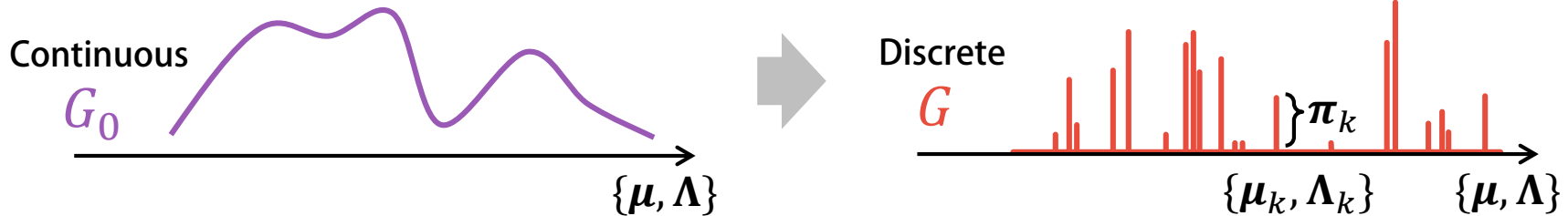$$\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k \sim G_0(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \quad \text{Gaussian-Wishart prior}$$

Continuous
$G_0$

$\{\boldsymbol{\mu}, \boldsymbol{\Lambda}\}$

Discrete
$G$

$\}\boldsymbol{\pi}_k$

$\{\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k\}$ $\{\boldsymbol{\mu}, \boldsymbol{\Lambda}\}$

- ## Generate samples independently

Equivalent representation

for $n = 1:N$
  $\boldsymbol{\mu}_n, \boldsymbol{\Lambda}_n \sim G(\boldsymbol{\mu}, \boldsymbol{\Lambda})$
  $\boldsymbol{x}_n \sim N(\boldsymbol{x}_n | \boldsymbol{\mu}_n, \boldsymbol{\Lambda}_n^{-1})$
end

$$\boldsymbol{z}_n \sim \mathrm{Categorical}(\boldsymbol{z}_n | \boldsymbol{\pi})$$

$$\boldsymbol{x}_n \sim \prod_{k=1}^{\infty} N(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)^{z_{nk}}$$

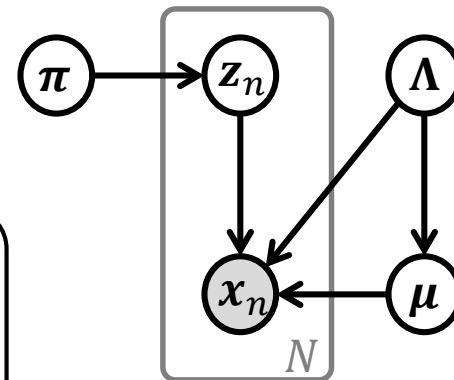Infinite GMM!

- Formulate a full joint distribution

$$p(\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = p(\boldsymbol{X}|\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})\,p(\boldsymbol{Z}|\boldsymbol{\pi})\,p(\boldsymbol{\pi})\,p(\boldsymbol{\mu}, \boldsymbol{\Lambda})$$

$$p(\boldsymbol{X}|\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{n=1}^{N}\prod_{k=1}^{K} N\!\left(\boldsymbol{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}\right)^{z_{nk}}$$

$$p(\boldsymbol{Z}|\boldsymbol{\pi}) = \prod_{n=1}^{N} \mathrm{Categorical}(\boldsymbol{z}_n|\boldsymbol{\pi}) = \prod_{n=1}^{N}\prod_{k=1}^{K} \pi_k^{z_{nk}}$$

Likelihood functions

$$p(\boldsymbol{\pi}) = \mathrm{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha}_0) = \frac{\Gamma\!\left(\sum_{k=1}^{K} \alpha_{0k}\right)}{\prod_{k=1}^{K} \Gamma(\alpha_{0k})} \prod_{k=1}^{K} \pi_k^{\alpha_{0k}-1}$$

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{k=1}^{K} N\!\left(\boldsymbol{\mu}_k|\boldsymbol{m}_0, (\beta_0 \boldsymbol{\Lambda}_k)^{-1}\right) W\!\left(\boldsymbol{\Lambda}_k|\boldsymbol{W}_0, \nu_0\right)$$

Prior distributions

- **Use a SBP prior instead of a Dirichlet prior**

$$p(X, Z, \pi, \mu, \Lambda, \alpha) = p(X|Z, \mu, \Lambda)\,p(Z|v)\,p(v|\alpha)\,p(\alpha)\,p(\mu, \Lambda)$$



$$p(X|Z, \mu, \Lambda) = \prod_{n=1}^{N} \prod_{k=1}^{\infty} N\left(x_n | \mu_k, \Lambda_k^{-1}\right)^{z_{nk}}$$

$$p(Z|v) = \prod_{n=1}^{N} \prod_{k=1}^{\infty} \left( \overbrace{v_k \prod_{k'=1}^{k-1} (1 - v_{k'})}^{\pi_k} \right)^{z_{nk}}$$

**Likelihood functions**

$$p(v|\alpha) = \prod_{k=1}^{\infty} \mathrm{Beta}(v_k | 1, \alpha) \quad p(\alpha) = \mathrm{Gamma}(\alpha | a_0, b_0)$$

**SBP prior**

$$p(\mu, \Lambda) = \prod_{k=1}^{\infty} N\left(\mu_k | m_0, (\beta_0 \Lambda_k)^{-1}\right) W(\Lambda_k | \mathbf{W}_0, \nu_0)$$

**Prior distributions**

- Beta-Bernoulli & Gamma-Exponential conjugacy
  - The VB is applicable for learning an iGMM
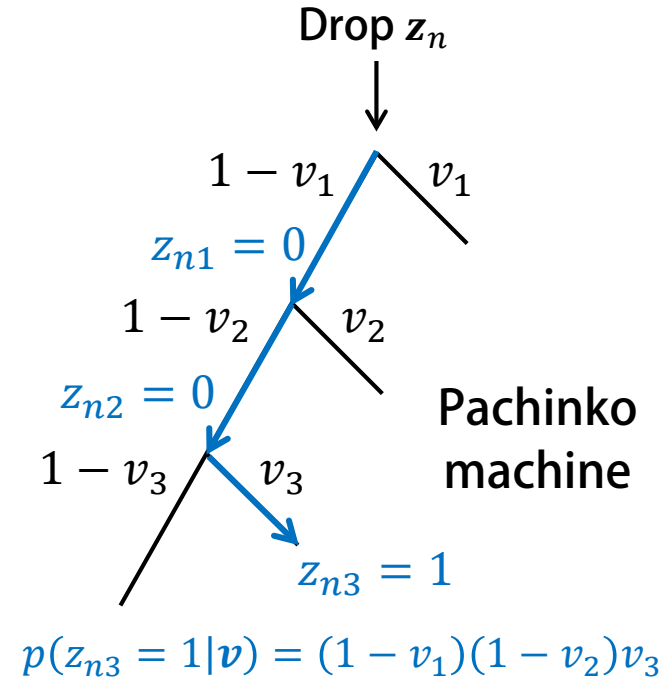
$$p(\mathbf{Z}|\mathbf{v}) = \prod_{n=1}^{N} \prod_{k=1}^{\infty} \left( v_k \prod_{k'=1}^{k-1} (1 - v_{k'}) \right)^{z_{nk}}$$

$$= \prod_{n=1}^{N} \prod_{k=1}^{\infty} v_k^{z_{nk}} (1 - v_k)^{\sum_{k'=k+1}^{\infty} z_{nk'}}$$

$$= \prod_{k=1}^{\infty} v_k^{\sum_{n=1}^{N} z_{nk}} (1 - v_k)^{\sum_{n=1}^{N} \sum_{k'>k} z_{nk'}}$$

The number of $z_n$'s that pass $k$

The number of $z_n$'s that stop at $k$

Conjugate

$$p(\mathbf{v}|\alpha) = \prod_{k=1}^{\infty} \alpha v_k^{1-1} (1 - v_k)^{\alpha - 1} \longleftrightarrow p(\alpha) = \frac{b_0^{a_0}}{\Gamma(a_0)} \alpha^{a_0 - 1} e^{-b_0 \alpha}$$

Conjugate

Drop $z_n$

$1 - v_1$    $v_1$

$z_{n1} = 0$

$1 - v_2$    $v_2$

$z_{n2} = 0$

Pachinko machine

$1 - v_3$    $v_3$

$z_{n3} = 1$

$$p(z_{n3} = 1|\mathbf{v}) = (1 - v_1)(1 - v_2)v_3$$

- Approximate a posterior $p(\boldsymbol{Z}, \boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda}, \alpha | \boldsymbol{X})$
  - Use a variational distribution $q(\boldsymbol{Z})q(\boldsymbol{v})q(\boldsymbol{\mu}, \boldsymbol{\Lambda})q(\alpha) \approx p(\boldsymbol{Z}, \boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda}, \alpha | \boldsymbol{X})$
  - Iteratively update (optimize) each factor
    - VB-E step
      - $\log q^*(\boldsymbol{Z}) = \langle \log p(\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda}, \alpha) \rangle_{q(\boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda}, \alpha)} + \text{const.}$
        $\qquad = \langle \log p(\boldsymbol{X} | \boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{Z} | \boldsymbol{v}) \rangle_{q(\boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} + \text{const.}$
    - VB-M step
      - $\log q^*(\boldsymbol{v}) = \langle \log p(\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}, \alpha)} + \text{const.}$
        $\qquad = \langle \log p(\boldsymbol{Z} | \boldsymbol{v}) p(\boldsymbol{v} | \alpha) \rangle_{q(\boldsymbol{Z}, \alpha)} + \text{const.}$
      - $\log q^*(\boldsymbol{\mu}, \boldsymbol{\Lambda}) = \langle \log p(\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{Z}, \boldsymbol{v}, \alpha)} + \text{const.}$
        $\qquad = \langle \log p(\boldsymbol{X} | \boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{Z})} + \text{const.}$
      - $\log q^*(\alpha) = \langle \log p(\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{Z}, \boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} + \text{const.}$
        $\qquad = \langle \log p(\boldsymbol{v} | \alpha) p(\alpha) \rangle_{q(\boldsymbol{v})} + \text{const.}$

- **Invoke the updating formula of VB**
  - Take the expectation of the full joint probability distribution under variational posteriors over other variables
  - Focus on only terms including $Z$
    (other terms can be absorbed into the normalization factor)

$$\log q^*(\boldsymbol{Z}) = \langle \log p(\boldsymbol{X}, \boldsymbol{Z}, \boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda}, \alpha) \rangle_{q(\boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda}, \alpha)} + \text{const.}$$
$$= \langle \log p(\boldsymbol{X}|\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{Z}|\boldsymbol{v}) p(\boldsymbol{v}|\boldsymbol{\alpha}) p(\alpha) p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda}, \alpha)} + \text{const.}$$
$$= \langle \log p(\boldsymbol{X}|\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{Z}|\boldsymbol{v}) \rangle_{q(\boldsymbol{v}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} + \text{const.}$$

$$p(\boldsymbol{X}|\boldsymbol{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{n=1}^{N} \prod_{k=1}^{\infty} N\left(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}\right)^{z_{nk}}$$

$$p(\boldsymbol{Z}|\boldsymbol{v}) = \prod_{n=1}^{N} \prod_{k=1}^{\infty} \left( v_k \prod_{k'=1}^{k-1} (1 - v_{k'}) \right)^{z_{nk}}$$

- Proceed the calculation according the updating rule

$$\langle \log p(\mathbf{Z}|\mathbf{v}) \rangle_{q(\mathbf{v})} = \sum_{n=1}^{N} \sum_{k=1}^{\infty} z_{nk} \left( \langle \log v_k \rangle_{q(v_k)} + \sum_{k'=1}^{k-1} \langle \log(1 - v_{k'}) \rangle_{q(v_{k'})} \right)$$

$$\langle \log p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{\mu},\boldsymbol{\Lambda})} = \sum_{n=1}^{N} \sum_{k=1}^{\infty} z_{nk} \left\langle \log N\left(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}\right) \right\rangle_{q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)}$$

$$\log q^*(\mathbf{Z}) = \langle \log p(\mathbf{Z}|\mathbf{v}) \rangle_{q(\mathbf{v})} + \langle \log p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \rangle_{q(\boldsymbol{\mu},\boldsymbol{\Lambda})} + \text{const.}$$

$$= \sum_{n=1}^{N} \sum_{k=1}^{\infty} z_{nk} \left( \underbrace{\langle \log v_k \rangle_{q(v_k)} + \sum_{k'=1}^{k-1} \langle \log(1 - v_{k'}) \rangle_{q(v_{k'})}}_{\text{Infinite GMM}} + \left\langle \log N\left(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}\right) \right\rangle_{q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)} \right) + \text{const.}$$

Infinite GMM

Finite GMM $\langle \log \pi_k \rangle_{q(\boldsymbol{\pi})}$

$$= \sum_{n=1}^{N} \sum_{k=1}^{\infty} z_{nk} \log \rho_{nk} + \text{const.}$$

- ## Calculate the variational posterior over latent variables $Z$
  - ▪ The normalization factor is automatically determined

$$\log q^*(\mathbf{Z}) = \sum_{n=1}^{N} \sum_{k=1}^{\infty} z_{nk} \log \rho_{nk} + \text{const.}$$

Truncate the variational posterior at the level $K$ *i.e.*, $q(z_{nk>K}) = 0$
The larger $K$ becomes, the more accurate the approximation is

$$\gamma_{nk} = \frac{\rho_{nk}}{\sum_{k'=1}^{K} \rho_{nk'}}$$

$$\log q^*(\mathbf{Z}) = \sum_{n=1}^{N} \sum_{k=1}^{\infty} z_{nk} \log \gamma_{nk}$$

$$q^*(\mathbf{Z}) = \prod_{n=1}^{N} \prod_{k=1}^{\infty} \gamma_{nk}^{z_{nk}} = \prod_{n=1}^{N} \text{Categorical}\,(\mathbf{z}_n | \boldsymbol{\gamma}_n)$$

Latent variables are categorical distributed!

- ## Invoke the updating formula of VB
  - Take the expectation of the full joint probability distribution under variational posteriors over other variables
  - Focus on only terms including $Z$
    (other terms can be absorbed into the normalization factor)

$$\log q^*(v) = \langle \log p(X, Z, v, \mu, \Lambda) \rangle_{q(Z,\mu,\Lambda,\alpha)} + \text{const.}$$
$$= \log p(v|\alpha) + \langle \log p(Z|v) \rangle_{q(Z)} + \text{const.}$$

$$\log q^*(\mu, \Lambda) = \langle \log p(X, Z, \pi, \mu, \Lambda) \rangle_{q(Z,\pi,\alpha)} + \text{const.} \quad \triangleleft \boxed{\text{Same as finite GMM}}$$
$$= \log p(\mu, \Lambda) + \langle \log p(X|Z, \mu, \Lambda) \rangle_{q(Z)} + \text{const.}$$

$$\log q^*(\alpha) = \langle \log p(X, Z, \pi, \mu, \Lambda) \rangle_{q(Z,v,\mu,\Lambda)} + \text{const.}$$
$$= \log p(\alpha) + \langle \log p(v|\alpha) \rangle_{q(v)} + \text{const.}$$

> **Bayesian estimation in simple conjugate models!**
> (Use responsibilities $q(Z)$ instead of latent variables $Z$)

- Calculate the variational posterior over parameters $\boldsymbol{v}$
  - The posteriors take the same forms of the priors

$$S_k[1] = \sum_{n=1}^{N} \gamma_{nk} \qquad S_k[\boldsymbol{x}] = \sum_{n=1}^{N} \gamma_{nk}\, \boldsymbol{x}_n \qquad S_k[\boldsymbol{x}\boldsymbol{x}^T] = \sum_{n=1}^{N} \gamma_{nk}\, \boldsymbol{x}_n \boldsymbol{x}_n^T$$

Sufficient statistics

$$p(\boldsymbol{v}|\alpha) = \prod_{k=1}^{\infty} \mathrm{Beta}(v_k|1, \alpha) = \prod_{k=1}^{\infty} \alpha v_k^{1-1} (1 - v_k)^{\alpha-1}$$

$$p(\boldsymbol{Z}|\boldsymbol{v}) = \prod_{k=1}^{\infty} v_k^{\sum_{n=1}^{N} z_{nk}} (1 - v_k)^{\sum_{n=1}^{N} \sum_{k'=k+1}^{\infty} z_{nk'}}$$

Bayes' theorem

$$p(\boldsymbol{v}|\boldsymbol{Z}, \alpha) = \prod_{k=1}^{\infty} \mathrm{Beta}\left(v_k \,\middle|\, 1 + \sum_{n=1}^{N} z_{nk}\,,\, \alpha + \sum_{n=1}^{N} \sum_{k'=k+1}^{\infty} z_{nk'}\right)$$

Replace $z_{nk}$ with $\gamma_{nk}$

$$q^*(\boldsymbol{v})$$

- Calculate the variational posterior over parameter $\alpha$
  - The posterior takes the same forms of the prior
  - Use that fact that if $x \sim \text{Beta}(1, \alpha)$, then $-\log(1-x) \sim \text{Exponential}(\alpha)$
  - $q^*(\alpha)$ is analytically tractable in case of iGMM

$$p(\alpha) = \text{Gamma}(\alpha|a_0, b_0) = \frac{b_0^{a_0}}{\Gamma(a_0)} \alpha^{a_0-1} e^{-b_0\alpha}$$

$$p(\boldsymbol{v}|\alpha) = \prod_{k=1}^{\infty} \text{Beta}(v_k|1, \alpha) = \alpha^K \prod_{k=1}^{\infty} (1-v_k)^{\alpha-1}$$

Bayes' theorem

$$p(\alpha|\boldsymbol{v}) = \text{Gamma}\left(\alpha\middle|a_0 + K, b_0 - \sum_{k=1}^{K} \log(1 - v_k)\right)$$

Replace $\log(1 - v_k)$ with $\langle \log(1 - v_k) \rangle_{q(v_k)}$

$$q^*(\alpha)$$

- **Truncate the variational poster $q(\boldsymbol{Z})$**
  - The infinite-dimensional true posterior $p(\boldsymbol{Z}|\boldsymbol{X})$ is NOT truncated!
  - $q(\boldsymbol{z}_n)$ is truncated at a sufficiently large level $K$ *i.e.,* $q(z_{nk>K}) = 0$
  - $K$ corresponds to how accurately $q(\boldsymbol{Z})$ approximates $p(\boldsymbol{Z}|\boldsymbol{X})$

- **Sort $K$ clusters in descending order before VB-M step**
  - Remove unnecessary cluster $k$ with $S_k[1] \approx 0$



This is effective for:
1. accelerating the convergence
2. avoiding poor local maxima

The SBP prior generates exponentially-decaying mixing ratios $\boldsymbol{\pi}$

- **Finite truncation at a certain level $K$ is required for VB**
  - A large amount of computational power is wasted
  - $K$ should be sufficiently large even if only a few clusters are required for representing the data

Infinite dim.

$v \sim \mathrm{GEM}(\alpha)$

$Z \sim \mathrm{CRP}(\alpha)$

If we can marginalize out $v$, we do not need to deal with infinity!

- **Marginalize out infinite-dimensional parameters $\pi$ or $v$**
  - Take the infinite limit of a Dirichlet-Categorical model

$K$-dimensional Dirichlet prior

$$\pi \sim \text{Dir}(\pi|\alpha\beta_K) \quad \beta_K = \underbrace{\left[\frac{1}{K}, \frac{1}{K}, \cdots, \frac{1}{K}\right]}_{K}$$

Likelihood

$$z_{1:N} \sim \text{Categorical}(z|\pi)$$

$\pi_3$   $\pi_2$

$\pi_1$

$\pi_5$   $\pi_4$

$\pi_6$

**Infinite-sided die**

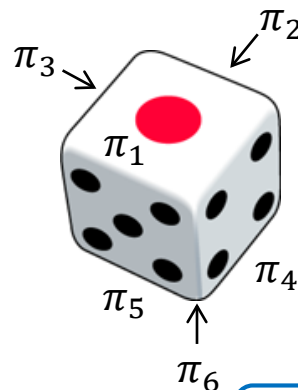Given $Z_{-n}$ as observed data, $z_n$ is predicted as:

$$p(z_{nk} = 1|Z_{-n}) = \int \underset{\text{Likelihood}}{p(z_{nk} = 1|\pi)} \underset{\text{Posterior}}{p(\pi|Z_{-n})} dZ_{-n}$$

The number of samples belonging to cluster $k$ among $N-1$ samples

$$= \int \pi_k \text{Dir}(\pi|\alpha\beta_K + \sum_{n' \neq n} z_{n'}) dZ_{-n} = \frac{\frac{\alpha}{K} + \overbrace{\sum_{n' \neq n} z_{n'k}}^{n_k^{(-n)}}}{\sum_{k'=1}^{K} \left(\frac{\alpha}{K} + \underbrace{\sum_{n' \neq n} z_{n'k'}}_{n_{k'}^{(-n)}}\right)} \quad \xrightarrow{K \to \infty} \quad \frac{n_k^{(-n)}}{(N-1) + \alpha}$$

- ## Focus on the probability that a new cluster is selected
  - ### Accumulate the probabilities that existing clusters are selected

$K$-dimensional Dirichlet prior

$$\pi \sim \mathrm{Dir}(\pi|\alpha\beta_K) \quad \beta_K = \underbrace{\left[\frac{1}{K}, \frac{1}{K}, \cdots, \frac{1}{K}\right]}_{K}$$

Likelihood

$$z_{1:N} \sim \mathrm{Categorical}(z|\pi)$$

$\pi_3$ $\pi_2$

$\pi_1$

$\pi_5$ $\pi_4$

$\pi_6$

→ **Infinite-sided die**

Given $Z_{-n}$ consisting of $K$ clusters, $z_n$ is predicted as:

Sum: $\frac{N-1}{N-1+\alpha}$

$$p(z_{nk} = 1|Z_{-n}) = \begin{cases} \dfrac{n_k^{(-n)}}{(N-1)+\alpha} & \text{Existing cluster } k \ (1 \leq k \leq K) \text{ is selected} \\[3mm] \dfrac{\alpha}{(N-1)+\alpha} & \text{New cluster } k \ (k > K) \text{ is created} \end{cases}$$

We index the new cluster as $K+1$

- Sequentially generate samples s.t. "the rich get richer"
  - Used as a prior on latent variables $Z$ ($= z_{1:N}$)

$$z_{1:N} \sim \mathrm{CRP}(\alpha) \quad \theta_k \sim G_0(\theta) \text{ if a new cluster is created}$$

Suppose $n - 1$ customers $z_{1:n-1}$ are already seated in restaurant $G_0$
The next customer $z_n$ stocastically selects a table as follows:

Existing tables

New table

7 customers

Dish $\theta_1$

3 customers

Dish $\theta_2$

5 customers

Dish $\theta_3$

Dish $\theta_4$

$G_0$

$p(z_{n1} = 1 | z_{1:n-1}) = \dfrac{7}{15 + \alpha}$ $\quad p(z_{n2} = 1 | z_{1:n-1}) = \dfrac{3}{15 + \alpha}$ $\quad p(z_{n3} = 1 | z_{1:n-1}) = \dfrac{5}{15 + \alpha}$ $\quad p(z_{n4} = 1 | z_{1:n-1}) = \dfrac{\alpha}{15 + \alpha}$

- The customer order does not change the CRP probability

$$\text{CRP}(\boldsymbol{Z}|\alpha) = p(\boldsymbol{z}_1)p(\boldsymbol{z}_2|\boldsymbol{z}_1)p(\boldsymbol{z}_3|\boldsymbol{z}_{1:2})p(\boldsymbol{z}_4|\boldsymbol{z}_{1:3})p(\boldsymbol{z}_5|\boldsymbol{z}_{1:4})p(\boldsymbol{z}_6|\boldsymbol{z}_{1:5})$$

3 customers          2 customers          1 customer

$p(z_{11} = 1) = \dfrac{\alpha}{0+\alpha}$ $\boldsymbol{z_1}$

$p(z_{42} = 1|\boldsymbol{z}_{1:3}) = \dfrac{\alpha}{3+\alpha}$ $\boldsymbol{z_4}$

$p(z_{63} = 1|\boldsymbol{z}_{1:5}) = \dfrac{\alpha}{5+\alpha}$ $\boldsymbol{z_6}$

Dish $\theta_1$   Dish $\theta_2$   Dish $\theta_3$

$\boldsymbol{z_2}$   $\boldsymbol{z_3}$   $\boldsymbol{z_5}$

$p(z_{21} = 1|\boldsymbol{z}_1) = \dfrac{1}{1+\alpha}$   $p(z_{31} = 1|\boldsymbol{z}_{1:2}) = \dfrac{2}{2+\alpha}$   $p(z_{52} = 1|\boldsymbol{z}_{1:4}) = \dfrac{1}{4+\alpha}$

$$\text{CRP}(\boldsymbol{Z}|\alpha) = p(\boldsymbol{z}_6)p(\boldsymbol{z}_5|\boldsymbol{z}_6)p(\boldsymbol{z}_4|\boldsymbol{z}_{5:6})p(\boldsymbol{z}_3|\boldsymbol{z}_{4:6})p(\boldsymbol{z}_2|\boldsymbol{z}_{3:6})p(\boldsymbol{z}_1|\boldsymbol{z}_{2:6})$$

3 customers          2 customers          1 customer

$p(z_{11} = 1|\boldsymbol{z}_{2:6}) = \dfrac{2+\alpha}{5+\alpha}$ $\boldsymbol{z_1}$

$p(z_{42} = 1|\boldsymbol{z}_{5:6}) = \dfrac{1+\alpha}{2+\alpha}$ $\boldsymbol{z_4}$

$p(z_{63} = 1) = \dfrac{\alpha}{0+\alpha}$ $\boldsymbol{z_6}$

Dish $\theta_1$   Dish $\theta_2$   Dish $\theta_3$

$\boldsymbol{z_2}$   $\boldsymbol{z_3}$   $\boldsymbol{z_5}$

$p(z_{21} = 1|\boldsymbol{z}_{3:6}) = \dfrac{1+\alpha}{4+\alpha}$   $p(z_{31} = 1|\boldsymbol{z}_{4:6}) = \dfrac{\alpha}{3+\alpha}$   $p(z_{52} = 1|\boldsymbol{z}_6) = \dfrac{\alpha}{1+\alpha}$

- Two major approaches to representing the DP
  - SBP: Represent how a distribution $G$ is drawn from the DP
  - CRP: Represent how samples $Z$ are drawn from the DP

Base measure
$G_0$

$G \sim \mathrm{DP}(\alpha, G_0)$

Infinite-dimensional distribution
$G$

$\theta_{1:N} \sim G$

Histogram
$Z$

$K$ clusters
$\theta$

$$G(\theta) = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}(\theta)$$

$$\pi_{1:\infty} \sim \mathrm{SBP}(\alpha) \quad \theta_{1:\infty} \sim G_0(\theta)$$

$$\mathbf{z}_{1:N} \sim \mathrm{CRP}(\alpha) \quad \theta_{1:K} \sim G_0(\theta)$$

- ## Reduce the number of variables for fast/better estimation
  - ### The parameters can be marginalized out because of conjugacy

$$p(X, Z, \mu, \Lambda, \alpha) = p(X|Z, \mu, \Lambda)p(Z|\alpha)p(\alpha)p(\mu, \Lambda) \Rightarrow p(X|Z) = p(X|Z)p(Z|\alpha)p(\alpha)$$

$p(Z|\alpha) = \mathrm{CRP}(Z|\alpha)$   Marginal likelihood for $Z$ (mixing ratios are marginalized out)

$$p(Z|\alpha) \propto \lim_{k \to \infty} \int p(Z|\pi)\mathrm{Dir}(\pi|\alpha\beta_K)d\pi$$

$p(\alpha) = \mathrm{Gamma}(\alpha|a_0, b_0)$   Hyper prior on $\alpha$

$$p(X|Z, \mu, \Lambda) = \prod_{n=1}^{N}\prod_{k=1}^{K} N(x_n|\mu_k, \Lambda_k^{-1})^{z_{nk}}$$

$$p(\mu, \Lambda) = \prod_{k=1}^{K} N(\mu_k|m_0, (\beta_0\Lambda_k)^{-1})W(\Lambda_k|\mathbf{W}_0, \nu_0)$$

Marginalization over $\mu, \Lambda$ is analytically tractable!

Conjugacy holds true
(Gaussian-Wishart-Gaussian)

- Generate samples from $p(\boldsymbol{Z}, \alpha | \boldsymbol{X})$
  - Divide $\{\boldsymbol{Z}, \alpha\}$ into $\{\boldsymbol{z}_1\}, \{\boldsymbol{z}_2\}, \cdots, \{\boldsymbol{z}_N\}, \{\alpha\}$
  - for $n = 1: N$
    - Sample $\boldsymbol{z}_n \sim p(\boldsymbol{z}_n | \boldsymbol{X}, \boldsymbol{Z}_{-n}, \alpha) = p(\boldsymbol{z}_n | \boldsymbol{x}_n, \boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n}, \alpha)$

$$p(z_{nk} = 1 | \boldsymbol{x}_n, \boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n}, \alpha) \propto p(z_{nk} = 1, \boldsymbol{x}_n | \boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n}, \alpha)$$

$$= p(z_{nk} = 1 | \boldsymbol{Z}_{-n}, \alpha) p(\boldsymbol{x}_n | z_{nk} = 1, \boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n})$$

$$= \mathrm{CRP}(z_{nk} = 1 | \boldsymbol{Z}_{-n}, \alpha) \int p(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) p(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k | \boldsymbol{X}_{-n}, \boldsymbol{Z}_{-n}) \, d\boldsymbol{\mu}_k d\boldsymbol{\Lambda}_k$$

$$= \begin{cases} \dfrac{n_k^{(-n)}}{N - 1 + \alpha} \mathrm{St}\left(\boldsymbol{x}_n \Big| \boldsymbol{m}_k^{(-n)}, \boldsymbol{L}_k^{(-n)}, v_k^{(-n)} + 1 - D\right) & \text{for existing cluster } k \ (1 \leq k \leq K) \\[3mm] \dfrac{\alpha}{N - 1 + \alpha} \mathrm{St}(\boldsymbol{x}_n | \boldsymbol{m}_0, \boldsymbol{L}_0, v_0 + 1 - D) & \text{for new cluster } K + 1 \end{cases}$$

- ## Update $z_n$ using the remove-and-add scheme
  - ▪ The number of tables $K$ **can be increased**



Seated at an existing table　　　　　　　Seated at a new table

- Update $z_n$ using the remove-and-add scheme
  - The number of tables $K$ can be decreased



$K = 2$

$z_1$   Dish $\theta_1$   $z_2$    $z_3$   Dish $\theta_2$

Remove $z_3$

$K = 1$

$z_1$   Dish $\theta_1$   $z_2$    $z_3$   Dish $\theta_2$

$\dfrac{2}{2 + \alpha}$    Add $z_3$    $\dfrac{\alpha}{2 + \alpha}$

**$K = 1$**

$z_1$   $z_3$   Dish $\theta_1$   $z_2$

Seated at an existing table

$K = 2$

$z_1$   Dish $\theta_1$   $z_2$    $z_3$   Dish $\theta_2^{new}$

Seated at a new table

- Calculate the probability of seating arrangement

$$p(\mathbf{Z}|\alpha) = \frac{1}{\sum_{i=1}^{N}(i-1+\alpha)} \prod_{k=1}^{K} \alpha(n_k - 1)! = \alpha^K \frac{\Gamma(\alpha)}{\Gamma(\alpha + N)} \prod_{k=1}^{K}(n_k - 1)!$$
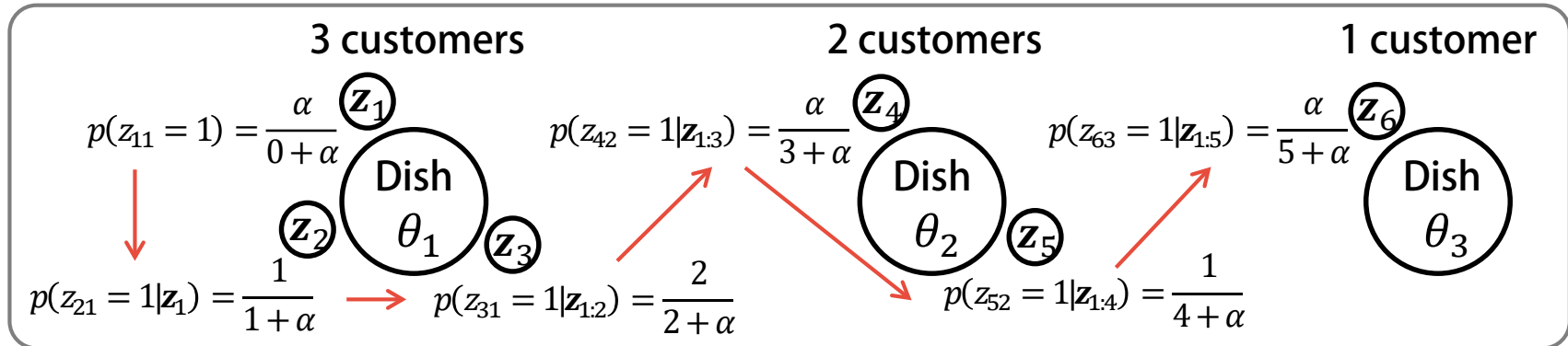
**Data augmentation**

$$p(\mathbf{Z}|\alpha) = \int p(\mathbf{Z}, \eta|\alpha)\, d\eta$$

$$p(\mathbf{Z}, \eta|\alpha) = \frac{\alpha^{K-1}(\alpha + n)}{\Gamma(N)} \eta^{\alpha}(1 - \eta)^{N-1} \prod_{k=1}^{K}(n_k - 1)!$$

$$1 = \int \text{Beta}(\eta|\alpha + 1, N)\, d\eta$$
$$= \frac{\Gamma(\alpha + N + 1)}{\Gamma(\alpha + 1)\Gamma(N)} \int \eta^{\alpha}(1 - \eta)^{N-1} d\eta$$
$$\Gamma(x + 1) = x\Gamma(x)$$



3 customers     2 customers     1 customer

$p(z_{11} = 1) = \frac{\alpha}{0 + \alpha}$   $z_1$

$p(z_{42} = 1|\mathbf{z}_{1:3}) = \frac{\alpha}{3 + \alpha}$   $z_4$

$p(z_{63} = 1|\mathbf{z}_{1:5}) = \frac{\alpha}{5 + \alpha}$   $z_6$

Dish $\theta_1$    Dish $\theta_2$    Dish $\theta_3$

$z_2$   $z_3$    $z_5$

$p(z_{21} = 1|\mathbf{z}_1) = \frac{1}{1 + \alpha}$   $p(z_{31} = 1|\mathbf{z}_{1:2}) = \frac{2}{2 + \alpha}$   $p(z_{52} = 1|\mathbf{z}_{1:4}) = \frac{1}{4 + \alpha}$

- Generate samples from $p(\mathbf{Z}, \alpha, \eta | \mathbf{X})$
  - Sample $\alpha \sim p(\alpha | \mathbf{X}, \mathbf{Z}, \eta) \propto p(\mathbf{Z}, \eta | \alpha) p(\alpha)$
  - Sample $\eta \sim p(\eta | \mathbf{X}, \mathbf{Z}, \alpha) \propto p(\mathbf{Z}, \eta | \alpha)$

**Sampling from beta**

$$p(\alpha) = \text{Gamma}(\alpha | a_0, b_0) = \frac{b_0^{a_0}}{\Gamma(a_0)} \alpha^{a_0 - 1} e^{-b_0 \alpha}$$

$$p(\eta | \mathbf{Z}, \alpha) = \text{Beta}(\alpha + 1, N)$$

$$p(\mathbf{Z}, \eta | \alpha) = \frac{\alpha^{K-1}(\alpha + n)}{\Gamma(N)} \eta^{\alpha} (1 - \eta)^{N-1} \prod_{k=1}^{K} (n_k - 1)!$$

**Bayes' theorem**

$$\propto \alpha^K \eta^{\alpha} + n \alpha^{K-1} \eta^{\alpha}$$

$$p(\alpha | \mathbf{Z}, \eta) \propto \alpha^{a_0 + K - 1} e^{-(b_0 - \log \eta)\alpha} + n \alpha^{a_0 + K - 2} e^{-(b_0 - \log \eta)\alpha}$$

$$\propto \omega \text{Gamma}(a_0 + K, b_0 - \log \eta) + (1 - \omega) \text{Gamma}(a_0 + K - 1, b_0 - \log \eta)$$

**Sampling from gamma mixture**

$$\frac{\omega}{1 - \omega} = \frac{a_0 + K - 1}{N(b_0 - \log \eta)}$$

- **Maximum likelihood estimation for finite GMM**
  - EM algorithm and hard EM (k-means)
- **Bayesian estimation for finite GMM**
  - (Collapsed) Gibbs sampling
  - (Collapsed) variational Bayes
- **Bayesian estimation for infinite GMM**
  - Collapsed Gibbs sampling with Chinese restaurant process
  - Variational Bayes with stick breaking process
- Other topics
  - Hierarchical Dirichlet process
    - HMM, PCFG (sequential data), LDA (grouped data)
  - Beta process, gamma process, Gaussian process
    - (Nonnegative) matrix factorization

> GS is feasible with SBP

> CVB is feasible with CRP

- Bayesian modeling
  - C. Bishop: Pattern Recognition and Machine Learning, Springer, 2010 (Ch. 9-11 & Appendix B).
  - 石井 健一郎, 上田 修功: 続・分かりやすいパターン認識 〜教師なし学習入門〜, オーム社, 2014.
  - 中島 伸一: 変分ベイズ学習, 講談社, 2016.
- Nonparametric Bayesian modeling
  - 佐藤 一誠: ノンパラメトリックベイズ 点過程と統計的機械学習の数理, 講談社, 2016.
  - D. Blei, M. Jordan: Variational Inference for Dirichlet Process Mixtures, Bayesian Analysis, Vol. 1, No. 1, pp.121-144, 2006.
  - J. Sung, Z. Ghahramani: Latent-Space Variational Bayes, IEEE Trans. on PAMI, Vol. 30, No. 12, 2008.
- Concentration parameter modeling
  - M. Escobar, M. West: Bayesian Density Estimation and Inference Using Mixtures, Journal of the American Statistical Association, Vol. 90, No. 430, pp. 577-588, 1995.
  - T. Stepleton: Understanding the Antoniak equation, 2008.
    http://www.cs.cmu.edu/~tss/antoniak.pdf

- ML estimation
  - Derive the update formulas of the parameters $\pi, \mu, \Lambda$ (p. 22) by letting the partial derivative of the lower bound (p. 20) w.r.t. each parameter equal to zero.
  - Implement the EM algorithm by using your favorite language.

- Bayesian estimation
  - Derive the variational posteriors of the parameters $\pi, \mu, \Lambda$ (p. 47) by using the formulas (p. 46)
  - Try one of the following at least:
    - Implement the VB algorithm
    - Implement the GS algorithm
  - Optional:
    - Implement the other algorithms for finite/infinite GMMs.

- **Report submission**
  - Deadline: 7/21 (Fri.)
  - "Assignments" → "Assignments 6/7 (Yoshii)"
  - Upload two files
    - PDF file: Report document
    - Zip file: Codes and instructions (README)
- **Program specification**
  - *your_program_or_script* x.csv z.csv params.dat
  - Show the value of the likelihood or lower bound at each iteration
  - Output z.csv and params.dat
    - z.csv: Posterior probabilities of $z_n$

```
0.2, 0.3, 0.5
0.5, 0.1, 0.4
0.1, 0.8, 0.1
    . . .
```