

BAYESIAN MIXTURE
MODELS

20

CHAPTER CONTENTS

Gaussian Mixture Models	221
Bayesian Formulation	221
Variational Inference	223
Gibbs Sampling	228
Latent Dirichlet Allocation (LDA)	229
Topic Models	230
Bayesian Formulation	231
Gibbs Sampling	232

In this chapter, within the Bayesian inference framework explained in [Chapter 17](#), practical Bayesian inference algorithms for mixture models are presented. First, algorithms based on variational approximation and Gibbs sampling for Gaussian mixture models are introduced. Then a variational Bayesian algorithm for topic models is explained.

20.1 GAUSSIAN MIXTURE MODELS

In this section, Bayesian inference algorithms for *Gaussian mixture models* (see [Section 15.1](#)) are introduced.

20.1.1 BAYESIAN FORMULATION

Let us consider the mixture of m Gaussian models:

$$q(\mathbf{x}|\mathcal{W}, \mathcal{M}, \mathcal{S}) = \sum_{\ell=1}^m w_{\ell} N(\mathbf{x}|\boldsymbol{\mu}_{\ell}, \mathbf{S}_{\ell}^{-1}),$$

where $\mathcal{W} = (w_1, \dots, w_m)$, $\mathcal{M} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_m)$, $\mathcal{S} = (\mathbf{S}_1, \dots, \mathbf{S}_m)$, and $N(\mathbf{x}|\boldsymbol{\mu}, \mathbf{S}^{-1})$ denotes the Gaussian density with expectation $\boldsymbol{\mu}$ and variance-covariance matrix \mathbf{S}^{-1} (or precision matrix \mathbf{S}):

$$N(\mathbf{x}|\boldsymbol{\mu}, \mathbf{S}^{-1}) = \frac{\sqrt{\det(\mathbf{S})}}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\top} \mathbf{S}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (20.1)$$

For training samples $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ drawn independently from the true density $p(\mathbf{x})$, the likelihood $p(\mathcal{D}|\mathcal{W}, \mathcal{M}, \mathcal{S})$ is given by

$$p(\mathcal{D}|\mathcal{W}, \mathcal{M}, \mathcal{S}) = \prod_{i=1}^n q(\mathbf{x}_i|\mathcal{W}, \mathcal{M}, \mathcal{S}).$$

For mixing weights \mathcal{W} , the *symmetric Dirichlet distribution* (see Section 6.3) is considered as the prior probability:

$$p(\mathcal{W}; \alpha_0) = \text{Dir}(\mathcal{W}; \alpha_0) \propto \prod_{\ell=1}^m w_{\ell}^{\alpha_0-1},$$

where $\text{Dir}(\mathcal{W}; \alpha)$ denotes the symmetric Dirichlet density with concentration parameter α . Note that the Dirichlet distribution is *conjugate* (see Section 17.2) for the discrete distribution given by Eq. (20.3).

For Gaussian expectations \mathcal{M} and Gaussian precision matrices \mathcal{S} , the product of the normal distribution and the Wishart distribution (see Section 6.4), called the *normal-Wishart distribution*, is considered as the prior probability:

$$\begin{aligned} p(\mathcal{M}, \mathcal{S}; \beta_0, \mathbf{W}_0, \nu_0) \\ &= \prod_{\ell=1}^m N(\boldsymbol{\mu}_{\ell} | \mathbf{0}, (\beta_0 \mathcal{S}_{\ell})^{-1}) W(\mathcal{S}_{\ell}; \mathbf{W}_0, \nu_0) \\ &\propto \prod_{\ell=1}^m \det(\mathcal{S}_{\ell})^{\frac{\nu_0-d}{2}-1} \exp\left(-\frac{\beta_0}{2} \boldsymbol{\mu}_{\ell}^{\top} \mathcal{S}_{\ell} \boldsymbol{\mu}_{\ell} - \frac{1}{2} \text{tr}(\mathbf{W}_0^{-1} \mathcal{S}_{\ell})\right), \end{aligned}$$

where $W(\mathcal{S}; \mathbf{W}, \nu)$ denotes the Wishart density with ν degrees of freedom:

$$W(\mathcal{S}; \mathbf{W}, \nu) = \frac{\det(\mathcal{S})^{\frac{\nu-d-1}{2}} \exp(-\frac{1}{2} \text{tr}(\mathbf{W}^{-1} \mathcal{S}))}{\det(2\mathbf{W})^{\frac{\nu}{2}} \Gamma_d(\frac{\nu}{2})}.$$

Here, d denotes the dimensionality of input \mathbf{x} and $\Gamma_d(\cdot)$ denotes the d -dimensional *gamma function* defined by Eq. (6.2):

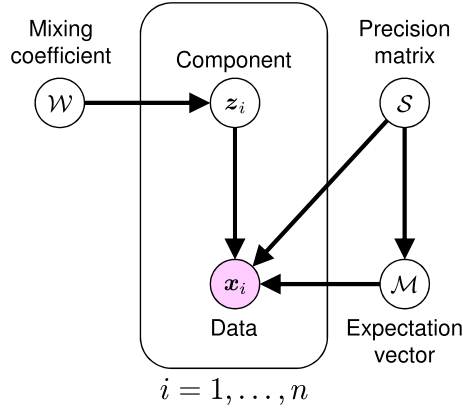
$$\Gamma_d\left(\frac{\nu}{2}\right) = \int_{\mathbb{S}_d^+} \det(\mathcal{S})^{\frac{\nu-d-1}{2}} \exp(-\text{tr}(\mathcal{S})) d\mathcal{S},$$

where \mathbb{S}_d^+ denotes the set of all $d \times d$ positive symmetric matrices. Note that the above *normal-Wishart distribution* is conjugate for the multivariate normal distribution with unknown expectation and unknown precision matrix.

The above formulation is summarized in Fig. 20.1. By the Bayes' theorem, the posterior probability $p(\mathcal{W}, \mathcal{M}, \mathcal{S}|\mathcal{D})$ is given as

$$\begin{aligned} p(\mathcal{W}, \mathcal{M}, \mathcal{S}|\mathcal{D}; \alpha_0, \beta_0, \mathbf{W}_0, \nu_0) \\ &= \frac{p(\mathcal{D}|\mathcal{W}, \mathcal{M}, \mathcal{S}) p(\mathcal{W}; \alpha_0) p(\mathcal{M}, \mathcal{S}; \beta_0, \mathbf{W}_0, \nu_0)}{p(\mathcal{D}; \alpha_0, \beta_0, \mathbf{W}_0, \nu_0)}, \end{aligned} \quad (20.2)$$

which is not computationally tractable. Below, practical approximate inference methods are introduced.

**FIGURE 20.1**

Variational Bayesian formulation of Gaussian mixture model.

To derive approximate inference methods, let us consider *latent variables*,

$$\mathcal{Z} = \{z_1, \dots, z_n\},$$

for training samples $\mathcal{D} = \{x_1, \dots, x_n\}$, where each $z_i = (z_{i,1}, \dots, z_{i,m})^\top$ is an m -dimensional vector that indicates mixture component selection. Specifically, if the k th component is selected, the k th element of z_i is 1 and all other elements are 0. The probability of observing \mathcal{Z} given mixing weights $\mathcal{W} = (w_1, \dots, w_m)$ can be expressed as

$$p(\mathcal{Z}|\mathcal{W}) = \prod_{i=1}^n \prod_{\ell=1}^m w_\ell^{z_{i,\ell}}. \quad (20.3)$$

20.1.2 VARIATIONAL INFERENCE

Here, the *variational technique* introduced in Section 18.2 is employed to approximately compute the posterior probability $p(\mathcal{W}, \mathcal{M}, \mathcal{S}|\mathcal{D})$ given by Eq. (20.2).

The marginal likelihood can be expressed as

$$\begin{aligned} \text{ML}(\alpha_0, \beta_0, \mathbf{W}_0, \nu_0) &= p(\mathcal{D}; \alpha_0, \beta_0, \mathbf{W}_0, \nu_0) \\ &= \iiint p(\mathcal{D}, \mathcal{Z}, \mathcal{W}, \mathcal{M}, \mathcal{S}; \alpha_0, \beta_0, \mathbf{W}_0, \nu_0) d\mathcal{Z} d\mathcal{W} d\mathcal{M} d\mathcal{S} \\ &= \iiint p(\mathcal{D}|\mathcal{Z}, \mathcal{M}, \mathcal{S}) p(\mathcal{Z}|\mathcal{W}) p(\mathcal{W}; \alpha_0) \\ &\quad \times p(\mathcal{M}, \mathcal{S}; \beta_0, \mathbf{W}_0, \nu_0) d\mathcal{Z} d\mathcal{W} d\mathcal{M} d\mathcal{S}. \end{aligned}$$

For the above marginal likelihood, let us consider the following trial distributions:

$$q(\mathcal{Z})r(\mathcal{W}, \mathcal{M}, \mathcal{S}).$$

Then, from Eq. (18.4), the VB-E step is given by

$$\begin{aligned} q(\mathcal{Z}) &\propto \exp \left(\iiint r(\mathcal{W}, \mathcal{M}, \mathcal{S}) \log p(\mathcal{D}, \mathcal{Z} | \mathcal{W}, \mathcal{M}, \mathcal{S}) d\mathcal{W} d\mathcal{M} d\mathcal{S} \right) \\ &= \exp \left(\iiint r(\mathcal{W}, \mathcal{M}, \mathcal{S}) \log p(\mathcal{D} | \mathcal{Z}, \mathcal{M}, \mathcal{S}) d\mathcal{W} d\mathcal{M} d\mathcal{S} \right. \\ &\quad \left. + \iiint r(\mathcal{W}, \mathcal{M}, \mathcal{S}) \log p(\mathcal{Z} | \mathcal{W}) d\mathcal{W} d\mathcal{M} d\mathcal{S} \right). \end{aligned} \quad (20.4)$$

Similarly, from Eq. (18.5), the VB-M step is given by

$$\begin{aligned} r(\mathcal{W}, \mathcal{M}, \mathcal{S}) &\propto p(\mathcal{W}, \mathcal{M}, \mathcal{S}; \alpha_0, \beta_0, \mathbf{W}_0, \nu_0) \\ &\quad \times \exp \left(\int q(\mathcal{Z}) \log p(\mathcal{D}, \mathcal{Z} | \mathcal{W}, \mathcal{M}, \mathcal{S}) d\mathcal{Z} \right) \\ &= p(\mathcal{W}; \alpha_0) \exp \left(\int q(\mathcal{Z}) \log p(\mathcal{Z} | \mathcal{W}) d\mathcal{Z} \right) \\ &\quad \times p(\mathcal{M}, \mathcal{S}; \beta_0, \mathbf{W}_0, \nu_0) \exp \left(\int q(\mathcal{Z}) \log p(\mathcal{D} | \mathcal{Z}, \mathcal{M}, \mathcal{S}) d\mathcal{Z} \right). \end{aligned} \quad (20.5)$$

Combining Eq. (20.4) and Eq. (20.5) yields the VBEM algorithm described in Fig. 20.2 (see Section 10.2.1 of reference [15] for details).

As mentioned in Section 18.2.1, $r(\mathcal{W}, \mathcal{M}, \mathcal{S})$ obtained by the VBEM algorithm may be a good approximation to the posterior probability $p(\mathcal{W}, \mathcal{M}, \mathcal{S} | \mathcal{D})$:

$$r(\mathcal{W}, \mathcal{M}, \mathcal{S}) = \text{Dir}(\mathcal{W} | \hat{\alpha}) \prod_{\ell=1}^m N(\boldsymbol{\mu}_\ell | \hat{\mathbf{h}}_\ell, (\hat{\beta}_\ell \mathbf{S}_\ell)^{-1}) W(\mathbf{S}_\ell | \hat{\mathbf{W}}_\ell, \hat{\nu}_\ell).$$

The expectations of the marginals of $r(\mathcal{W}, \mathcal{M}, \mathcal{S})$ can be obtained as

$$\hat{w}_\ell = \frac{\hat{\alpha}_\ell}{\sum_{\ell'=1}^m \hat{\alpha}_{\ell'}}, \quad \hat{\boldsymbol{\mu}}_\ell = \hat{\mathbf{h}}_\ell, \quad \text{and} \quad \hat{\mathbf{S}}_\ell = \hat{\nu}_\ell \hat{\mathbf{W}}_\ell,$$

which may be used as the most plausible parameter values. Then the following density estimator is obtained:

$$\hat{p}(\mathbf{x}) = \sum_{\ell=1}^m \hat{w}_\ell N(\mathbf{x} | \hat{\boldsymbol{\mu}}_\ell, \hat{\mathbf{S}}_\ell^{-1}).$$

A MATLAB code for computing this VBEM-based density estimator is given in Fig. 20.3, and its behavior is illustrated in Fig. 20.4. Here, the mixture model of five Gaussian components is fitted to the mixture of two Gaussian distributions. As shown in Fig. 20.4, two out of five Gaussian components fit the true two Gaussian distributions well, and the remaining three Gaussian components are almost eliminated—the learned mixing coefficients are given as

1. Initialize parameters $\{\hat{\alpha}_\ell, \hat{\mathbf{h}}_\ell, \hat{\beta}_\ell, \hat{\mathbf{W}}_\ell, \hat{\nu}_\ell\}_{\ell=1}^m$.
2. VB-E step: Compute the distribution of $\mathcal{Z} = \{z_1, \dots, z_n\}$ from current solution $\{\hat{\alpha}_\ell, \hat{\mathbf{h}}_\ell, \hat{\beta}_\ell, \hat{\mathbf{W}}_\ell, \hat{\nu}_\ell\}_{\ell=1}^m$:

$$q(\mathcal{Z}) = \prod_{i=1}^n \prod_{\ell=1}^m \hat{\eta}_{i,\ell}^{z_{i,\ell}}, \quad \text{where} \quad \hat{\eta}_{i,\ell} = \frac{\hat{\rho}_{i,\ell}}{\sum_{\ell'=1}^m \hat{\rho}_{i,\ell'}}.$$

$\{\hat{\eta}_{i,\ell}\}_{i=1}^n, \ell=1}^m$ are the responsibilities computed as

$$\begin{aligned} \hat{\rho}_{i,\ell} = \exp & \left(\psi(\hat{\alpha}_\ell) - \psi\left(\sum_{\ell'=1}^m \hat{\alpha}_{\ell'}\right) + \frac{1}{2} \sum_{j=1}^d \psi\left(\frac{\hat{\nu}_\ell + 1 - j}{2}\right) \right. \\ & \left. + \frac{1}{2} \log \det(\hat{\mathbf{W}}_\ell) - \frac{d}{2\hat{\beta}_\ell} - \frac{\hat{\nu}_\ell}{2} (\mathbf{x}_i - \hat{\mathbf{h}}_\ell)^\top \hat{\mathbf{W}}_\ell (\mathbf{x}_i - \hat{\mathbf{h}}_\ell) \right), \end{aligned}$$

where $\psi(\alpha)$ denotes the *digamma function* defined as the log-derivative of the gamma function $\Gamma(\alpha)$:

$$\psi(\alpha) = \frac{d}{d\alpha} \log \Gamma(\alpha) = \frac{\Gamma'(\alpha)}{\Gamma(\alpha)}.$$

3. VB-M step: Compute the joint distribution of $\mathcal{W} = (w_1, \dots, w_m)$, $\mathcal{M} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_m)$, and $\mathcal{S} = (S_1, \dots, S_m)$ from current responsibilities $\{\hat{\eta}_{i,\ell}\}_{i=1}^n, \ell=1}^m$:

$$r(\mathcal{W}, \mathcal{M}, \mathcal{S}) = \text{Dir}(\mathcal{W} | \hat{\boldsymbol{\alpha}}) \prod_{\ell=1}^m N(\boldsymbol{\mu}_\ell | \hat{\mathbf{h}}_\ell, (\hat{\beta}_\ell S_\ell)^{-1}) W(S_\ell | \hat{\mathbf{W}}_\ell, \hat{\nu}_\ell),$$

where

$$\begin{aligned} \hat{\gamma}_\ell &= \sum_{i=1}^n \hat{\eta}_{i,\ell}, \quad \hat{\mathbf{c}}_\ell = \frac{1}{\hat{\gamma}_\ell} \sum_{i=1}^n \hat{\eta}_{i,\ell} \mathbf{x}_i, \quad \hat{\mathbf{h}}_\ell = \frac{\hat{\gamma}_\ell}{\hat{\beta}_\ell} \hat{\mathbf{c}}_\ell, \\ \hat{\alpha}_\ell &= \alpha_0 + \hat{\gamma}_\ell, \quad \hat{\beta}_\ell = \beta_0 + \hat{\gamma}_\ell, \quad \hat{\nu}_\ell = \nu_0 + \hat{\gamma}_\ell, \quad \text{and} \\ \hat{\mathbf{W}}_\ell &= \left(\mathbf{W}_0^{-1} + \sum_{i=1}^n \hat{\eta}_{i,\ell} (\mathbf{x}_i - \hat{\mathbf{c}}_\ell)(\mathbf{x}_i - \hat{\mathbf{c}}_\ell)^\top + \frac{\beta_0 \hat{\gamma}_\ell}{\beta_0 + \hat{\gamma}_\ell} \hat{\mathbf{c}}_\ell \hat{\mathbf{c}}_\ell^\top \right)^{-1}. \end{aligned}$$

4. Iterate 2–3 until convergence.

FIGURE 20.2

VBEM algorithm for Gaussian mixture model. $(\alpha_0, \beta_0, \mathbf{W}_0, \nu_0)$ are hyperparameters.

```

x=[2*randn(1,100)-5 randn(1,50); randn(1,100) randn(1,50)+3];
[d,n]=size(x); m=5; e=rand(n,m); W=zeros(d,d,m); b0=1;
for o=1:10000
    e=e./repmat(sum(e,2),[1 m]);
    g=sum(e); a=1+g; b=b0+g; nu=3+g; w=a/sum(a);
    xe=x*e; c=xe./repmat(g,[d 1]); h=xe./repmat(b,[d 1]);
    for k=1:m
        t1=x-repmat(c(:,k),[1 n]); t2=x-repmat(h(:,k),[1 n]);
        W(:, :, k)=inv(eye(d)+(t1.*repmat(e(:,k)',[d 1]))*t1' ...
            +c(:,k)*c(:,k)'*b0*g(k)/(b0+g(k)));
        t3=sum(psi((nu(k)+1-[1:d])/2))+log(det(W(:, :, k)));
        e(:,k)=exp(t3/2+psi(a(k))-psi(sum(a))-d/2/b(k) ...
            -sum(t2.*(W(:, :, k)*t2))*nu(k)/2);
    end
    if o>1 && norm(w-w0)+norm(h-h0)+norm(W(:)-W0(:))<0.001
        break
    end
    w0=w; h0=h; W0=W;
end

figure(1); clf; hold on
plot(x(1,:),x(2,:), 'ro'); v=linspace(0,2*pi,100);
for k=1:m
    [V,D]=eig(nu(k)*W(:, :, k));
    X=3*w(k)*V'*[cos(v)/D(1,1); sin(v)/D(2,2)];
    plot(h(1,k)+X(1,:),h(2,k)+X(2,:), 'b-')
end

```

FIGURE 20.3

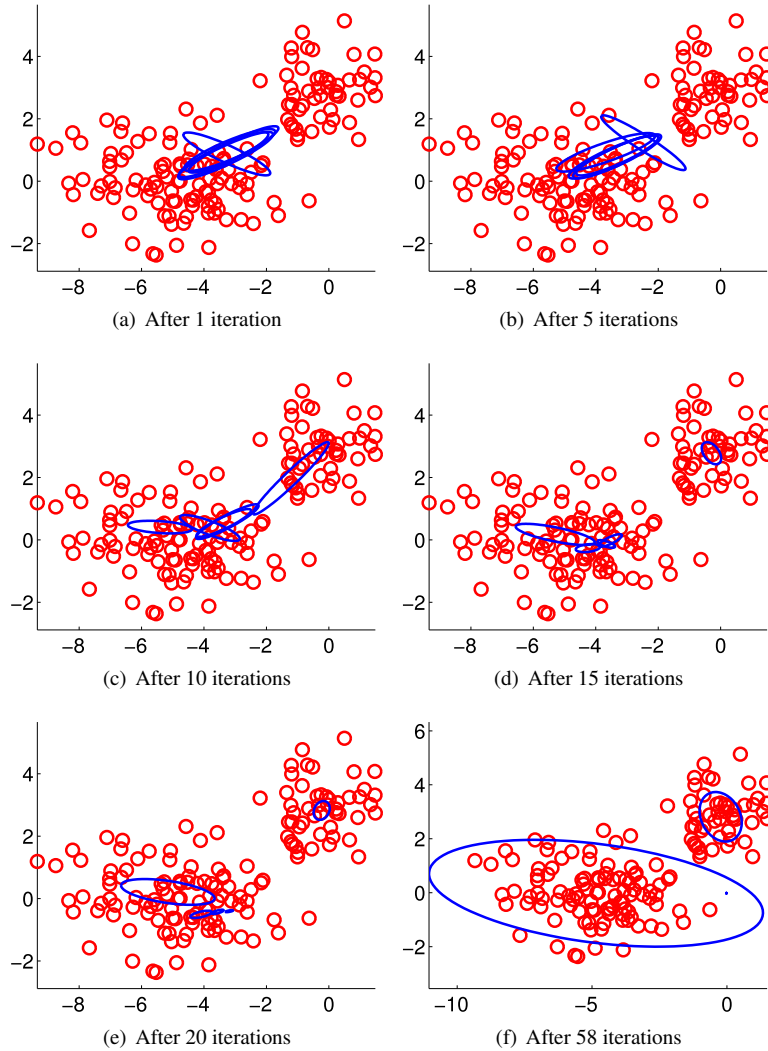
MATLAB code of VBEM algorithm for Gaussian mixture model.

$$(\hat{w}_1, \hat{w}_2, \hat{w}_3, \hat{w}_4, \hat{w}_5) = (0.01, 0.33, 0.01, 0.01, 0.65).$$

As explained in Section 18.2.1, the negative of the variational free energy,

$$\begin{aligned}
 F(q, r) = & \iiint q(\mathcal{Z}) r(\mathcal{W}, \mathcal{M}, \mathcal{S}) \\
 & \times \log \frac{q(\mathcal{Z}) r(\mathcal{W}, \mathcal{M}, \mathcal{S})}{p(\mathcal{D}, \mathcal{Z}, \mathcal{W}, \mathcal{M}, \mathcal{S}; \alpha_0, \beta_0, \mathbf{W}_0, \nu_0)} d\mathcal{Z} d\mathcal{W} d\mathcal{M} d\mathcal{S},
 \end{aligned}$$

gives a lower bound of the log marginal likelihood, $\text{ML}(\alpha_0, \beta_0, \mathbf{W}_0, \nu_0)$. Based on the dependency in Fig. 20.1, the variational free energy can be expressed as

**FIGURE 20.4**

Example of VBEM algorithm for Gaussian mixture model. The size of ellipses is proportional to the mixing weights $\{w_\ell\}_{\ell=1}^m$. A mixture model of five Gaussian components is used here, but three components have mixing coefficient close to zero and thus they are almost eliminated.

$$\begin{aligned}
 F(q, r) = & \iiint q(\mathcal{Z}) r(\mathcal{W}, \mathcal{M}, \mathcal{S}) \left(\log q(\mathcal{Z}) + \log r(\mathcal{W}, \mathcal{M}, \mathcal{S}) \right. \\
 & - \log p(\mathcal{D} | \mathcal{Z}, \mathcal{M}, \mathcal{S}) - \log p(\mathcal{Z} | \mathcal{W}) - \log p(\mathcal{W}; \alpha_0) \\
 & \left. - \log p(\mathcal{M}, \mathcal{S}; \beta_0, \mathbf{W}_0, \nu_0) \right) d\mathcal{Z} d\mathcal{W} d\mathcal{M} d\mathcal{S},
 \end{aligned}$$

which allows us to choose the hyperparameters and the number of mixing components by the empirical Bayes method (see Section 17.4). See Section 10.2.2 of reference [15] for more details of computing the variational free energy for Gaussian mixture models.

20.1.3 GIBBS SAMPLING

Next, the *Gibbs sampling technique* introduced in Section 19.3.3 is employed to numerically approximate the posterior probability $p(\mathcal{W}, \mathcal{M}, \mathcal{S} | \mathcal{D})$ given by Eq. (20.2).

Collapsed Gibbs sampling for assignment z_i can be performed as follows [74]:

$$p(z_{i,\ell} = 1 | \tilde{\mathcal{Z}}, \mathcal{D}; \alpha_0, \beta_0, \mathbf{W}_0, \nu_0) \propto (\hat{\alpha}_\ell - 1) \times t\left(\mathbf{x}_i | \hat{\mathbf{c}}_\ell, \frac{\hat{\beta}_\ell + 1}{\hat{\beta}_\ell(\hat{\nu}_\ell - d + 1)} \hat{\mathbf{W}}_\ell^{-1}, \hat{\nu}_\ell - d + 1\right), \quad (20.6)$$

where $\tilde{\mathcal{Z}} = \{z_{i'}\}_{i' \neq i}$, $\hat{\alpha}_\ell = \alpha_0 + n_\ell$, $\hat{\beta}_\ell = \beta_0 + n_\ell$, $\hat{\nu}_\ell = \nu_0 + n_\ell$,

$$\hat{\mathbf{W}}_\ell^{-1} = \left(\mathbf{W}_0 + \sum_{i: z_{i,\ell}=1} \mathbf{x}_i \mathbf{x}_i^\top - \frac{n_\ell^2}{\hat{\beta}_\ell} \hat{\mathbf{c}}_\ell \hat{\mathbf{c}}_\ell^\top \right)^{-1},$$

$n_\ell = \sum_{i=1}^n z_{i,\ell}$, $\hat{\mathbf{c}}_\ell = \frac{1}{n_\ell} \sum_{i: z_{i,\ell}=1} \mathbf{x}_i$, and $t(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu)$ denotes the probability density function of the *multivariate Student's t-distribution*:

$$t(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu) = \frac{\Gamma\left(\frac{\nu+d}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)(\nu\pi)^{d/2} \det(\boldsymbol{\Sigma})^{1/2}} \left(1 + \frac{1}{\nu} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)^{-\frac{\nu+d}{2}}.$$

After estimating assignment \mathcal{Z} , parameters $\mathcal{W}, \mathcal{M}, \mathcal{S}$ for each Gaussian component may be estimated separately as

$$\hat{\mathbf{w}}_\ell = \frac{\hat{\alpha}_\ell}{\sum_{\ell'=1}^m \hat{\alpha}_{\ell'}}, \quad \hat{\boldsymbol{\mu}}_\ell = \frac{n_\ell}{\hat{\beta}_\ell} \hat{\mathbf{c}}_\ell, \quad \text{and} \quad \hat{\mathbf{S}}_\ell = \hat{\nu}_\ell \hat{\mathbf{W}}_\ell.$$

A MATLAB code of collapsed Gibbs sampling for Gaussian mixture models is given in Fig. 20.5, and its behavior is illustrated in Fig. 20.6. A mixture model of five Gaussian components is used here, but only two components remain and no samples belong to the remaining three components after Gibbs sampling.

The *Dirichlet process*, denoted by $\text{DP}(\alpha_0, p_0)$, is a probability distribution of probability distributions, and a Dirichlet distribution is generated from a Dirichlet process. p_0 is called the *base distribution* which corresponds to the mean of the Dirichlet distribution, while α_0 is called the *concentration parameter* which corresponds to the *precision* (i.e., the inverse variance) of the Dirichlet distribution. The *Dirichlet process mixture* is a mixture model using the Dirichlet process as a prior probability and is equivalent to considering an infinite-dimensional Dirichlet distribution. Indeed, taking the limit $m \rightarrow \infty$ for $\alpha_0 = \alpha'_0/m$ in Eq. (20.6) gives a *collapsed Gibbs sampling* procedure for the Dirichlet process mixture model.


```

x=[2*randn(1,100)-5 randn(1,50); randn(1,100) randn(1,50)+3];
[d,n]=size(x); m=5; z=mod(randperm(n),m)+1;
a0=1; b0=1; n0=1; W0=eye(d);
for o=1:100
    for i=1:n
        g=(1:n~i); X=x(:,g); Z=z(g);
        for k=1:m
            p(k)=0; e=(Z==k); t=sum(e);
            if t~=0
                u=n0+t-d+1; b=b0+t; c=sum(X(:,e),2); xi=x(:,i)-c/t;
                W=inv((b+1)/b/u*(W0+X(:,e)*X(:,e)'+c*c'/b));
                p(k)=(a0+t-1)*gamma((u+d)/2)/gamma(u/2)*u^(-d/2) ...
                    *sqrt(det(W))*(1+xi'*W*xi/u)^(-(u+d)/2);
            end
        end
        z(i)=find(cumsum(p/sum(p))>rand,1);
    end
end

figure(1); clf; hold on
plot(x(1,:),x(2,:), 'ro'); v=linspace(0,2*pi,100);
for k=1:m
    e=(z==k); t=sum(e); nu(k)=n0+t; u=nu(k)-d+1; b=b0+t;
    c=sum(x(:,e),2); w(k)=a0+t; h(:,k)=c/b; W(:, :, k)=zeros(d);
    if t~=0
        W(:, :, k)=inv((W0+x(:,e)*x(:,e)'+c*c'/b));
    end
end
w=w./sum(w);
for k=1:m
    [V,D]=eig(nu(k)*W(:, :, k));
    X=3*w(k)*V'*[cos(v)/D(1,1); sin(v)/D(2,2)];
    plot(h(1,k)+X(1,:),h(2,k)+X(2,:), 'b-')
end

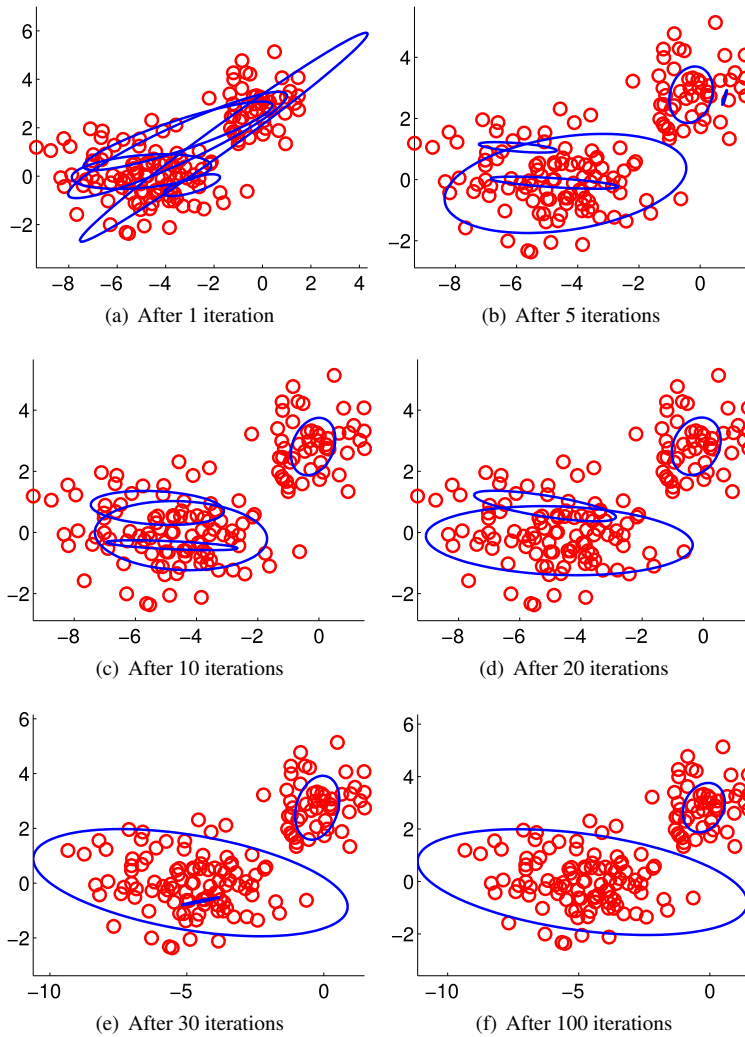
```

FIGURE 20.5

MATLAB code of collapsed Gibbs sampling for Gaussian mixture model.

20.2 LATENT DIRICHLET ALLOCATION (LDA)

One of the most successful applications of the Bayesian generative approach is *topic modeling* called LDA [16]. In this section, an implementation of LDA based on *Gibbs sampling* is explained [51].

**FIGURE 20.6**

Example of collapsed Gibbs sampling for Gaussian mixture model. A mixture model of five Gaussian components is used here, but only two components remain and no samples belong to the remaining three components.

20.2.1 TOPIC MODELS

A topic model is a generative model for a set of D documents,

$$\mathcal{W} = \{w_1, \dots, w_D\},$$

where the d th document $\mathbf{w}_d = (w_{d,1}, \dots, w_{d,n_d})$ is a sequence of n_d words in a corpus with W unique words. More specifically, each word $w_{n,d} \in \{1, \dots, W\}$ is supposed to be given by the mixture of T topics as

$$p(w_{n,d} | \boldsymbol{\theta}^{(d)}, \boldsymbol{\phi}^{(1)}, \dots, \boldsymbol{\phi}^{(T)}) = \sum_{t_{n,d}=1}^T p(t_{n,d} | \boldsymbol{\theta}^{(d)}) p(w_{n,d} | t_{n,d}, \boldsymbol{\phi}^{(t_{n,d})}),$$

where $t_{n,d} \in \{1, \dots, T\}$ is a latent variable that indicates the topic of $w_{n,d}$. For topic $t = 1, \dots, T$, the probability,

$$p(t | \boldsymbol{\theta}^{(d)}) = \theta_t^{(d)} \geq 0,$$

denotes the discrete probability of topic t with parameter $\boldsymbol{\theta}^{(d)} = (\theta_1^{(d)}, \dots, \theta_T^{(d)})^\top$ such that $\sum_{t=1}^T \theta_t^{(d)} = 1$. For $w = 1, \dots, W$, the probability,

$$p(w | t, \boldsymbol{\phi}^{(t)}) = \phi_w^{(t)} \geq 0,$$

denotes the discrete probability of observing word w under topic t with parameter $\boldsymbol{\phi}^{(t)} = (\phi_1^{(t)}, \dots, \phi_W^{(t)})^\top$ such that $\sum_{w=1}^W \phi_w^{(t)} = 1$. This generative model means that, for each word w in document d , topic t is chosen following the discrete probability with parameter $\boldsymbol{\theta}^{(d)}$ and then word w is chosen following the discrete probability with parameter $\boldsymbol{\phi}^{(t)}$.

The goal of topic modeling is to specify the multinomial parameters $\boldsymbol{\theta}^{(d)}$ and $\boldsymbol{\phi}^{(t)}$ from data \mathcal{W} .

20.2.2 BAYESIAN FORMULATION

Here, instead of directly learning the parameters $\boldsymbol{\theta}^{(d)}$ and $\boldsymbol{\phi}^{(t)}$, let us consider the posterior probability of the set \mathcal{T} of all topics for all documents \mathcal{W} :

$$\mathcal{T} = \{\mathbf{t}_1, \dots, \mathbf{t}_D\},$$

where $\mathbf{t}_d = (t_{d,1}, \dots, t_{d,n_d})$.

For the multinomial parameters $\boldsymbol{\theta}^{(d)}$ and $\boldsymbol{\phi}^{(t)}$, the *symmetric Dirichlet distributions* (see Section 6.3) are considered as the prior probabilities:

$$p(\boldsymbol{\theta}^{(d)}; \boldsymbol{\alpha}) = \text{Dir}(\boldsymbol{\theta}^{(d)}; \boldsymbol{\alpha}) \quad \text{and} \quad p(\boldsymbol{\phi}^{(t)}; \boldsymbol{\beta}) = \text{Dir}(\boldsymbol{\phi}^{(t)}; \boldsymbol{\beta}),$$

where $\text{Dir}(\cdot; \boldsymbol{\alpha})$ denotes the symmetric Dirichlet density with concentration parameter $\boldsymbol{\alpha}$. Note that the Dirichlet distributions are *conjugate* (see Section 17.2) for discrete distributions.

The posterior probability of topics \mathcal{T} is given by

$$p(\mathcal{T} | \mathcal{W}; \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{p(\mathcal{W}, \mathcal{T}; \boldsymbol{\alpha}, \boldsymbol{\beta})}{\sum_{\mathcal{T}} p(\mathcal{W}, \mathcal{T}; \boldsymbol{\alpha}, \boldsymbol{\beta})}, \quad (20.7)$$

where the marginal likelihood $p(\mathcal{W}, \mathcal{T}; \alpha, \beta)$ can be expressed *analytically* as

$$\begin{aligned} p(\mathcal{W}, \mathcal{T}; \alpha, \beta) &= \prod_{d=1}^D \prod_{n=1}^{n_d} \int \theta_{t_{n,d}}^{(d)} \text{Dir}(\theta^{(d)}; \alpha) d\theta^{(d)} \int \phi_{w_{n,d}}^{(t_{n,d})} \text{Dir}(\phi^{(t_{n,d})}; \beta) d\phi^{(t_{n,d})} \\ &= \left(\frac{\Gamma(T\alpha)}{\Gamma(\alpha)^T} \right)^D \prod_{d=1}^D \frac{\prod_{t=1}^T \Gamma(n_t^{(d)} + \alpha)}{\Gamma(\sum_{t=1}^T n_t^{(d)} + T\alpha)} \left(\frac{\Gamma(W\beta)}{\Gamma(\beta)^W} \right)^T \prod_{t=1}^T \frac{\prod_{w=1}^W \Gamma(n_t^{(w)} + \beta)}{\Gamma(\sum_{w=1}^W n_t^{(w)} + W\beta)}, \end{aligned}$$

where $n_t^{(d)}$ denotes the number of times a word from document d is assigned to topic t , $n_t^{(w)}$ denotes the number of times word w is assigned to topic t , and $\Gamma(\cdot)$ denotes the *gamma function* (see Section 4.3).

20.2.3 GIBBS SAMPLING

Let us use *Gibbs sampling* (see Section 19.3.3) to approximate the posterior probability $p(\mathcal{T}|\mathcal{W}; \alpha, \beta)$.

Let $\tilde{\mathcal{W}}$ and $\tilde{\mathcal{T}}$ be defined in the same way as \mathcal{W} and \mathcal{T} , but $w_{n,d}$ and $t_{n,d}$ are excluded. Then the conditional probability needed for Gibbs sampling is given by

$$\begin{aligned} p(t_{n,d}|\tilde{\mathcal{T}}, \mathcal{W}) &= p(t_{n,d}|\tilde{\mathcal{T}}, w_{n,d}, \tilde{\mathcal{W}}) \\ &\propto p(t_{n,d}|\tilde{\mathcal{T}}) p(w_{n,d}|t_{n,d}, \tilde{\mathcal{T}}, \tilde{\mathcal{W}}). \end{aligned} \quad (20.8)$$

The first term in Eq. (20.8) can be expressed as

$$p(t_{n,d}|\tilde{\mathcal{T}}) = \int p(t_{n,d}|\theta^{(d)}) p(\theta^{(d)}|\tilde{\mathcal{T}}) d\theta^{(d)}, \quad (20.9)$$

where Bayes' theorem yields

$$p(\theta^{(d)}|\tilde{\mathcal{T}}) \propto p(\tilde{\mathcal{T}}|\theta^{(d)}) p(\theta^{(d)}; \alpha).$$

Here, the prior probability $p(\theta^{(d)}; \alpha)$ is chosen to be symmetric Dirichlet distribution with concentration parameter α , which is conjugate for the discrete distribution $p(\tilde{\mathcal{T}}|\theta^{(d)})$. Then the posterior probability $p(\theta^{(d)}|\tilde{\mathcal{T}})$ is also the symmetric Dirichlet distribution with concentration parameter $\tilde{n}_{t_{n,d}}^{(d)} + \alpha$, where $\tilde{n}_{t_{n,d}}^{(d)}$ is defined in the same way as $n_{t_{n,d}}^{(d)}$, but $w_{n,d}$ and $t_{n,d}$ are excluded. Then Eq. (20.9) can be computed analytically as

$$p(t_{n,d}|\tilde{\mathcal{T}}) = \int \theta_{t_{n,d}}^{(d)} \text{Dir}(\theta^{(d)}; \tilde{n}_{t_{n,d}}^{(d)} + \alpha) d\theta^{(d)} = \frac{\tilde{n}_{t_{n,d}}^{(d)} + \alpha}{\sum_{t=1}^T \tilde{n}_t^{(d)} + T\alpha}.$$

Similarly, the second term in Eq. (20.8) can be expressed as

$$p(w_{n,d}|t_{n,d}, \tilde{\mathcal{T}}, \tilde{\mathcal{W}}) = \int p(w_{n,d}|t_{n,d}, \phi^{(t_{n,d})}) p(\phi^{(t_{n,d})}|\tilde{\mathcal{T}}, \tilde{\mathcal{W}}) d\phi^{(t_{n,d})}, \quad (20.10)$$

where Bayes' theorem yields

$$p(\phi^{(t_{n,d})} | \tilde{\mathcal{T}}, \tilde{\mathcal{W}}) \propto p(\tilde{\mathcal{W}} | \phi^{(t_{n,d})}, \tilde{\mathcal{T}}) p(\phi^{(t_{n,d})}; \beta).$$

Here, the prior probability $p(\phi^{(t_{n,d})}; \beta)$ is chosen to be the symmetric Dirichlet distribution with concentration parameter β , which is conjugate for the discrete distribution $p(\tilde{\mathcal{W}} | \phi^{(t_{n,d})}, \tilde{\mathcal{T}})$. Then the posterior probability $p(\phi^{(t_{n,d})} | \tilde{\mathcal{T}}, \tilde{\mathcal{W}})$ is also the symmetric Dirichlet distribution with concentration parameter $\tilde{n}_{t_{n,d}}^{(w_{n,d})} + \beta$, where $\tilde{n}_{t_{n,d}}^{(w_{n,d})}$ is defined in the same way as $n_{t_{n,d}}^{(w_{n,d})}$, but $w_{n,d}$ and $t_{n,d}$ are excluded. Then Eq. (20.10) can be computed analytically as

$$\begin{aligned} p(w_{n,d} | t_{n,d}, \tilde{\mathcal{T}}, \tilde{\mathcal{W}}) &= \int \phi_{w_{n,d}}^{(t_{n,d})} \text{Dir}(\phi^{(t_{n,d})}; \tilde{n}_{t_{n,d}}^{(w_{n,d})} + \beta) d\phi^{(t_{n,d})} \\ &= \frac{\tilde{n}_{t_{n,d}}^{(w_{n,d})} + \beta}{\sum_{w=1}^W \tilde{n}_{t_{n,d}}^{(w)} + W\beta}. \end{aligned}$$

Summarizing the above equations, Eq. (20.8) can be expressed as

$$p(t_{n,d} | \tilde{\mathcal{T}}, \mathcal{W}) \propto \frac{\tilde{n}_{t_{n,d}}^{(d)} + \alpha}{\sum_{t=1}^T \tilde{n}_t^{(d)} + T\alpha} \times \frac{\tilde{n}_{t_{n,d}}^{(w_{n,d})} + \beta}{\sum_{w=1}^W \tilde{n}_{t_{n,d}}^{(w)} + W\beta},$$

which shows that Gibbs sampling can be efficiently performed just by counting the number of words.

Finally, with estimated topics $\hat{\mathcal{T}}$, the solutions $\hat{\theta}_t^{(d)}$ and $\hat{\phi}_t^{(w)}$ for any $w \in \{1, \dots, W\}$, $t \in \{1, \dots, T\}$, and $d \in \{1, \dots, D\}$ can be computed as

$$\hat{\theta}_t^{(d)} = \frac{\tilde{n}_{t_{n,d}}^{(d)} + \alpha}{\sum_{t=1}^T \tilde{n}_t^{(d)} + T\alpha} \quad \text{and} \quad \hat{\phi}_t^{(w)} = \frac{\tilde{n}_{t_{n,d}}^{(w_{n,d})} + \beta}{\sum_{w=1}^W \tilde{n}_{t_{n,d}}^{(w)} + W\beta}.$$

