

GAUSSIAN MIXTURE MODEL FOR CLUSTERED OBSERVATIONS

ANIKO SZABO

1. MIXTURE MODEL

Let Y_{ij} , $i = 1, \dots, \mathcal{I}$, $j = 1, \dots, n_i$ denote a continuous random variable measured in subject j of cluster i . We want to model the distribution of Y as a mixture of K latent normal distributions:

$$[Y_{ij} \mid \xi_{ij} = k] \sim N(\mu_k, \sigma_k), \quad (1)$$

where ξ_{ij} is the latent class indicator for subject j of cluster i taking values from 1 to K .

Due to the clustering we cannot assume that ξ_{ij} are independent, however we will assume that they follow and exchangeable multinomial distribution $\mathcal{EM}(K, \mathbf{q})$.

Denoting $\phi(y; \mu, \sigma)$ the normal pdf with mean μ and standard deviation σ , the likelihood is

$$L = \prod_{i=1}^{\mathcal{I}} \sum_{\mathbf{x} \in \mathbb{Z}_K^{n_i}} Pr(\boldsymbol{\xi}_i = \mathbf{x}) \times \prod_{j=1}^{n_i} \phi(y_{ij}; \mu_{x_j}, \sigma_{x_j}) \quad (2)$$

$$= \prod_{i=1}^{\mathcal{I}} \sum_{\mathbf{x} \in \mathbb{Z}_K^{n_i}} \frac{q_{\dot{\mathbf{x}}|n_i}}{\binom{n_i}{\dot{\mathbf{x}}}} \times \prod_{j=1}^{n_i} \phi(y_{ij}; \mu_{x_j}, \sigma_{x_j}), \quad (3)$$

where $\mathbf{x} = (x_1, \dots, x_{n_i})$, $\boldsymbol{\xi}_i = (\xi_{i1}, \dots, \xi_{in_i})$, $\dot{\mathbf{x}} = (\sum_j I(x_j = 1), \dots, \sum_j I(x_j = K))$ is the vector of latent category frequencies, and

$$q_{\mathbf{r}|n_i} = \sum_{\mathbf{t} \in \mathbb{Z}_N: \mathbf{t}'\mathbf{1} = N} h(\mathbf{r}, \mathbf{t}, N, n_i) q_{\mathbf{t}|L} \quad (4)$$

are the exchangeable and marginally reproducible outcome probabilities of the exchangeable multinomial distribution with maximal cluster size N with $h(\mathbf{r}, \mathbf{t}, N, n_i)$ defined as multivariate hypergeometric probabilities.

2. EM ALGORITHM

We can consider ξ_{ij} missing data. Then the complete data likelihood is

$$L_c(\mu, \sigma, q | \xi_{ij}, y_{ij}) = \sum_{i=1}^{\mathcal{I}} \frac{q_{\boldsymbol{\xi}|n_i}}{\binom{n_i}{\boldsymbol{\xi}}} \prod_{j=1}^{n_i} \log \phi(y_{ij}; \mu_{\xi_{ij}}, \sigma_{\xi_{ij}}) = \sum_{i=1}^{\mathcal{I}} \prod_{\mathbf{z} \in \mathbb{Z}_K^{n_i}} \left[\frac{q_{\dot{\mathbf{z}}|n_i}}{\binom{n_i}{\dot{\mathbf{z}}}} \prod_{j=1}^{n_i} \log \phi(y_{ij}; \mu_{z_j}, \sigma_{z_j}) \right]^{I(\boldsymbol{\xi}_i = \mathbf{z})}. \quad (5)$$

Its logarithm without the binomial term that does not contain unknown parameters is

$$\log L_c(\mu, \sigma, q | \xi_{ij}, y_{ij}) = \sum_{i=1}^{\mathcal{I}} \sum_{\mathbf{z} \in \mathbb{Z}_K^{n_i}} I(\boldsymbol{\xi}_i = \mathbf{z}) \left[\log q_{\dot{\mathbf{z}}|n_i} + \sum_{j=1}^{n_i} \log \phi(y_{ij}; \mu_{z_j}, \sigma_{z_j}) \right] \quad (6)$$

The expected complete data log-likelihood, given previous estimates $\mu^{(m)}, \sigma^{(m)}, q^{(m)}$ is

$$Q(\mu, \sigma, q) = E[\log L_c(\mu, \sigma, q \mid \mu^{(m)}, \sigma^{(m)}, q^{(m)}, y_{ij})] = \sum_{i=1}^{\mathcal{I}} \sum_{\mathbf{z} \in \mathbb{Z}_K^{n_i}} Pr(\boldsymbol{\xi}_i = \mathbf{z} \mid \mu^{(m)}, \sigma^{(m)}, q^{(m)}, y_{ij}) [\log q_{\dot{\mathbf{z}}|n_i} + \sum_{j=1}^{n_i} \log \phi(y_{ij}; \mu_{z_j}, \sigma_{z_j})], \quad (7)$$

2.1. **E-step.** Using the Bayes theorem

$$e_{iz}^{(m)} = Pr(\xi_i = z \mid \mu^{(m)}, \sigma^{(m)}, q^{(m)}, y_{ij}) = \frac{Pr(y_i \mid \xi_i = z, \mu^{(m)}, \sigma^{(m)}, q^{(m)}) Pr(\xi_i = z \mid \mu^{(m)}, \sigma^{(m)}, q^{(m)})}{\sum_w Pr(y_i \mid \xi_i = w, \mu^{(m)}, \sigma^{(m)}, q^{(m)}) Pr(\xi_i = w \mid \mu^{(m)}, \sigma^{(m)}, q^{(m)})} = \frac{\prod_{j=1}^{n_i} \phi(y_{ij}; \mu_{z_j}^{(m)}, \sigma_{z_j}^{(m)}) q_{z|n_i}^{(m)} / \binom{n_i}{z}}{\sum_w \prod_{j=1}^{n_i} \phi(y_{ij}; \mu_{w_j}^{(m)}, \sigma_{w_j}^{(m)}) q_{w|n_i}^{(m)} / \binom{n_i}{w}}. \quad (8)$$

2.2. **M-step.**

$$Q(\mu, \sigma, q) = \sum_{i=1}^{\mathcal{I}} \sum_{z \in \mathbb{Z}_K^{n_i}} e_{iz}^{(m)} [\log q_{z|n_i} + \sum_{j=1}^{n_i} \log \phi(y_{ij}; \mu_{z_j}, \sigma_{z_j})] = \quad (10)$$

$$\sum_{i=1}^{\mathcal{I}} \sum_{z \in \mathbb{Z}_K^{n_i}} e_{iz}^{(m)} \log q_{z|n_i} + \sum_{i=1}^{\mathcal{I}} \sum_{z \in \mathbb{Z}_K^{n_i}} e_{iz}^{(m)} \sum_{j=1}^{n_i} \log \phi(y_{ij}; \mu_{z_j}, \sigma_{z_j}) = \quad (11)$$

$$\sum_{i=1}^{\mathcal{I}} \sum_{z \in \mathbb{Z}_K^{n_i}} e_{iz}^{(m)} \log q_{z|n_i} + \sum_{k=1}^K \sum_{i=1}^{\mathcal{I}} \sum_{j=1}^{n_i} \sum_{z: z_j=k} e_{iz}^{(m)} \log \phi(y_{ij}; \mu_k, \sigma_k), \quad (12)$$

where the two components can be maximized separately, and within the last component the terms corresponding to different k 's can be maximized separately as well.

2.2.1. *Update for q .* In terms of q , the log-likelihood can be viewed as the \mathcal{EM} log-likelihood for an extended data set: for each cluster i and for each possible frequency vector of responses $\mathbf{r} \in \mathbb{Z}_{n_i}$, $\mathbf{r}'\mathbf{1} = n_i$ we compute the 'observed' frequency as $a_{i\mathbf{r}}^{(m)} = \sum_{z: \mathbf{z}=\mathbf{r}} e_{iz}^{(m)}$. Then $q_{\mathbf{t}|N}^{(m+1)}$ can be obtained from the EM algorithm for fitting the \mathcal{EM} model to these frequencies.

2.2.2. *Update for μ and σ .* In terms of the normal distribution parameters, we have a weighted normal log-likelihood for each k . The weight corresponding to observation j in cluster i is $b_{ijk}^{(m)} = \sum_{z: z_j=k} e_{iz}^{(m)}$. Then

$$\mu_k^{(m+1)} = \sum_{i=1}^{\mathcal{I}} \sum_{j=1}^{n_i} b_{ijk}^{(m)} y_{ij} \quad (13)$$

$$[\sigma_k^{(m+1)}]^2 = \sum_{i=1}^{\mathcal{I}} \sum_{j=1}^{n_i} b_{ijk}^{(m)} (y_{ij} - \mu_k^{(m+1)})^2 \quad (14)$$