

Homework 4

Instructor: Ke Tran

Email: m.k.tran@uva.nl

Student: Andrea Jemmett

UvA ID: 11162929

Collaborators: N/A

Email: andrea.jemmett@student.uva.nl

You are allowed to discuss with your colleagues but you should write the answers in *your own words*. If you discuss with others, write down the name of your collaborators on top of the first page. No points will be deducted for collaborations. If we find similarities in solutions beyond the listed collaborations we will consider it as cheating. We will not accept any late submissions under any circumstances. The solutions to the previous homework will be handed out in the class at the beginning of the next homework session. After this point, late submissions will be automatically graded zero.

Problem 1. Consider a Gaussian mixture model

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

1. Given the expected value of the complete-data log-likelihood (9.40 in Bishop's book)

$$\mathbb{E}_{\text{posterior}}[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\}$$

Derive update rules for $\boldsymbol{\pi}$, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$.

2. Consider a special case of the model above, in which the covariance matrices $\boldsymbol{\Sigma}_k$ of the components are all constrained to have a common value $\boldsymbol{\Sigma}$. Derive EM equations for maximizing the likelihood function under such a model.

Problem 2. Suppose we wish to use the EM algorithm to maximize the posterior distribution

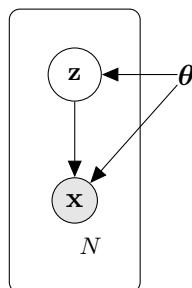


Figure 1: A simple generative model.

$p(\boldsymbol{\theta} | \mathbf{X})$ for a model (Figure 1) containing latent variables \mathbf{z} and observed variables \mathbf{x} . Show that the

E step remains the same as in the maximum likelihood case, where as in the M step, the quantity to be maximized is

$$\sum_{\mathbf{z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta})$$

Problem 3.

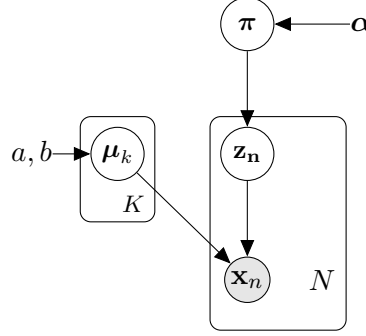


Figure 2: Mixtures of Bernoulli distribution

$$\begin{aligned} \boldsymbol{\pi}|\boldsymbol{\alpha} &\sim \text{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha}) \\ \mathbf{z}_n|\boldsymbol{\pi} &\sim \text{Mult}(\mathbf{z}_n|\boldsymbol{\pi}) \\ \boldsymbol{\mu}_k|a_k, b_k &\sim \text{Beta}(\boldsymbol{\mu}_k|a_k, b_k) \\ \mathbf{x}_n|\mathbf{z}_n, \boldsymbol{\mu} = \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K\} &\sim \prod_{k=1}^K (\text{Bern}(\mathbf{x}_n|\boldsymbol{\mu}_k))^{z_{nk}} \end{aligned}$$

Derive the EM algorithm for maximizing the posterior probability $p(\boldsymbol{\mu}, \boldsymbol{\pi}|\{\mathbf{x}_n\}_{n=1}^N)$. (The E step is given in Bishop's Book, you only need to do the M step)

Solution: The M step consists in the maximization of the expectation of the complete log-likelihood. In the case of the model in Figure 2 it is given by

$$\begin{aligned} \mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\pi})] &= \ln \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} + \sum_{k=1}^K \left[D \ln \frac{\Gamma(a_k + b_k)}{\Gamma(a_k)\Gamma(b_k)} + (\alpha_k - 1) \ln \pi_k \right. \\ &\quad \left. + (a_k - 1) \sum_{i=1}^D \ln \mu_{ki} + (b_k - 1) \sum_{i=1}^D (1 - x_{ni}) \ln(1 - \mu_{ki}) \right] \\ &\quad + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left[\ln \pi_k + \sum_{i=1}^D x_{ni} \ln \mu_{ki} + \sum_{i=1}^D (1 - x_{ni}) \ln(1 - \mu_{ki}) \right] \end{aligned} \quad (1)$$

In the M step we need to maximize the expectation of the complete data log-likelihood w.r.t. the

parameters μ_k and π . To do this we evaluate the partial derivatives w.r.t. μ_{ki}

$$\begin{aligned} \frac{\partial}{\partial \mu_{ki}} \mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}, \mu, \pi)] &= \frac{a_k - 1}{\mu_{ki}} - \frac{b_k - 1}{1 - \mu_{ki}} + \sum_{n=1}^N \sum_{k=1}^K \left\{ \gamma(z_{nk}) \right. \\ &\quad \left. \times \frac{\partial}{\partial \mu_{ki}} \left[\ln \pi_k + \sum_{i=1}^D x_{ni} \ln \mu_{ki} + \sum_{i=1}^D (1 - x_{ni}) \ln(1 - \mu_{ki}) \right] \right\} \end{aligned} \quad (2)$$

$$= \frac{a_k - 1}{\mu_{ki}} - \frac{b_k - 1}{1 - \mu_{ki}} + \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{x_{ni}}{\mu_{ki}} - \frac{1 - x_{ni}}{1 - \mu_{ki}} \right] \quad (3)$$

$$= \frac{a_k - 1}{\mu_{ki}} - \frac{b_k - 1}{1 - \mu_{ki}} + \frac{\sum_{n=1}^N \gamma(z_{nk}) x_{ni}}{\mu_{ki}} - \frac{\sum_{n=1}^N (1 - x_{ni}) \gamma(z_{nk})}{1 - \mu_{ki}} \quad (4)$$

$$= \frac{1}{\mu_{ki}(1 - \mu_{ki})} \left[(1 - \mu_{ki})(a_k - 1) - \mu_{ki}(b_k - 1) + (1 - \mu_{ki}) \sum_{n=1}^N \gamma(z_{nk}) x_{ni} - \mu_{ki} \sum_{n=1}^N \gamma(z_{nk}) (1 - x_{ni}) \right] \quad (5)$$

$$= \frac{1}{\mu_{ki}(1 - \mu_{ki})} \left[a_k - 1 - \mu_{ki} a_k + \mu_{ki} - \mu_{ki} b_k + \mu_{ki} + \sum_{n=1}^N \gamma(z_{nk}) x_{ni} - \mu_{ki} \sum_{n=1}^N \gamma(z_{nk}) x_{ni} - \mu_{ki} \sum_{n=1}^N \gamma(z_{nk}) + \mu_{ki} \sum_{n=1}^N \gamma(z_{nk}) x_{ni} \right] \quad (6)$$

and by setting this derivative to zero and solving for μ_{ki}

$$0 = a_k - 1 - \mu_{ki} a_k + \mu_{ki} - \mu_{ki} b_k + \mu_{ki} + \sum_{n=1}^N \gamma(z_{nk}) x_{ni} - \mu_{ki} \sum_{n=1}^N \gamma(z_{nk}) \quad (7)$$

$$1 - a_k - \sum_{n=1}^N \gamma(z_{nk}) x_{ni} = \mu_{ki} \left[2 - a_k - b_k - \sum_{n=1}^N \gamma(z_{nk}) \right] \quad (8)$$

$$\mu_{ki} = \frac{1 - a_k - \sum_{n=1}^N \gamma(z_{nk}) x_{ni}}{2 - a_k - b_k - \sum_{n=1}^N \gamma(z_{nk})} = \frac{\sum_{n=1}^N \gamma(z_{nk}) x_{ni} + a_k - 1}{\sum_{n=1}^N \gamma(z_{nk}) + a_k + b_k - 2} \quad (9)$$

which rewritten in vector form gives

$$\mu_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n + a_k - 1}{\sum_{n=1}^N \gamma(z_{nk}) + a_k + b_k - 2} \quad (10)$$

Next we need to maximize the expectation of the log-likelihood w.r.t. π_k using Lagrange multipliers to enforce the constraint $\sum_{k=1}^K \pi_k = 1$

$$\frac{\partial}{\partial \pi_k} \mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}, \mu, \pi)] = \frac{\alpha_k - 1}{\pi_k} + \frac{\sum_{n=1}^N \gamma(z_{nk})}{\pi_k} + \lambda \quad (11)$$

$$= \frac{\alpha_k - 1 + \sum_{n=1}^N \gamma(z_{nk}) + \lambda \pi_k}{\pi_k} \quad (12)$$

and by setting this derivative to zero and solving for π_k

$$0 = \alpha_k - 1 + \sum_{n=1}^N \gamma(z_{nk}) + \lambda \pi_k \quad (13)$$

$$\pi_k = \frac{1 - \alpha_k - \sum_{n=1}^N \gamma(z_{nk})}{\lambda} \quad (14)$$

Using the constraint on $\boldsymbol{\pi}$ we can find the value of λ

$$\sum_{k=1}^K \pi_k = 1 \quad (15)$$

$$\frac{1}{\lambda} \sum_{k=1}^K \left[1 - \alpha_k - \sum_{n=1}^N \gamma(z_{nk}) \right] = 1 \quad (16)$$

$$\lambda = K - \sum_{k=1}^K \alpha_k - \sum_{k=1}^K \sum_{n=1}^N \gamma(z_{nk}) \quad (17)$$

$$\lambda = K - \sum_{k=1}^K \alpha_k - N \quad (18)$$

substituting back we get

$$\pi_k = \frac{1 - \alpha_k - \sum_{n=1}^N \gamma(z_{nk})}{K - \sum_{k=1}^K \alpha_k - N} \quad (19)$$

■

Problem 4★. Bishop 10.38 Verify the results of the calculation of mean, variance of $q^{\setminus n}(\boldsymbol{\theta})$ and normalizing constant Z_n for the expectation propagation algorithm applied to the cluster problem

$$\mathbf{m}^{\setminus n} = \mathbf{m} + v^{\setminus n} v_n^{-1} (\mathbf{m} - \mathbf{m}_n) \quad (20)$$

$$(v^{\setminus n})^{-1} = v^{-1} - v_n^{-1} \quad (21)$$

$$Z_n = (1 - w) \mathcal{N}(\mathbf{x}_n | \mathbf{m}^{\setminus n}, (v^{\setminus n} + 1) \mathbf{I}) + w \mathcal{N}(\mathbf{x}_n | \mathbf{0}, a \mathbf{I}) \quad (22)$$

Hint: see 10.38 in the book for more hint. You need equation (2.115) in Bishop to solve this homework.

Solution: The first step is to use the division formula (eq. 10.205 in Bishop) to derive $q^{\setminus n}(\boldsymbol{\theta})$

$$q^{\setminus n}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{\tilde{f}_n(\boldsymbol{\theta})} \quad (23)$$

$$= \frac{\mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, v\mathbf{I})}{s_n \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}_n, v_n\mathbf{I})} \quad (24)$$

$$\propto \exp \left\{ \frac{1}{2v_n}(\boldsymbol{\theta} - \mathbf{m}_n)^\top (\boldsymbol{\theta} - \mathbf{m}_n) - \frac{1}{2v}(\boldsymbol{\theta} - \mathbf{m})^\top (\boldsymbol{\theta} - \mathbf{m}) \right\} \quad (25)$$

$$= \exp \left\{ \frac{1}{2v_n}(\boldsymbol{\theta}^\top \boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \mathbf{m}_n + \mathbf{m}_n^\top \mathbf{m}_n) - \frac{1}{2v}(\boldsymbol{\theta}^\top \boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \mathbf{m} + \mathbf{m}^\top \mathbf{m}) \right\} \quad (26)$$

$$= \exp \left\{ \frac{1}{2} \left(\frac{1}{v_n} - \frac{1}{v} \right) \boldsymbol{\theta}^\top \boldsymbol{\theta} - \boldsymbol{\theta}^\top \left(\frac{\mathbf{m}_n}{v_n} - \frac{\mathbf{m}}{v} \right) + \frac{1}{2v_n} \mathbf{m}_n^\top \mathbf{m}_n - \frac{1}{2v} \mathbf{m}^\top \mathbf{m} \right\} \quad (27)$$

$$= \exp \left\{ -\frac{1}{2}(v^{-1} - v_n^{-1})\boldsymbol{\theta}^\top \boldsymbol{\theta} + \boldsymbol{\theta}^\top (v^{-1}\mathbf{m} - v_n^{-1}\mathbf{m}_n) + \text{const} \right\} \quad (28)$$

$$\propto \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}^{\setminus n}, v^{\setminus n}\mathbf{I}) \quad (29)$$

which is a Gaussian distribution with mean and variance given by

$$(v^{\setminus n})^{-1} = v^{-1} - v_n^{-1} \quad (30)$$

$$(v^{\setminus n})^{-1}\mathbf{m}^{\setminus n} = v^{-1}\mathbf{m} - v_n^{-1}\mathbf{m}_n \quad (31)$$

$$\mathbf{m}^{\setminus n} = (v^{-1}\mathbf{m} - v_n^{-1}\mathbf{m}_n)v^{\setminus n} \quad (32)$$

To evaluate Z_n we can use eq. 10.206 in Bishop

$$Z_n = \int q^{\setminus n}(\boldsymbol{\theta}) f_n(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (33)$$

$$= \int \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}^{\setminus n}, v^{\setminus n}\mathbf{I}) \left[(1-w)\mathcal{N}(\mathbf{x}_n|\boldsymbol{\theta}, \mathbf{I}) + w\mathcal{N}(\mathbf{x}_n|\mathbf{0}, a\mathbf{I}) \right] d\boldsymbol{\theta} \quad (34)$$

We note that $q^{\setminus n}(\boldsymbol{\theta})$ is a Gaussian distribution over $\boldsymbol{\theta}$ and $f_n(\boldsymbol{\theta})$ is a Gaussian distribution over \mathbf{x}_n conditioned on $\boldsymbol{\theta}$. We can then interpret Z_n in the following way

$$Z_n = \int p(\boldsymbol{\theta}) p(\mathbf{x}_n|\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (35)$$

$$= \int p(\boldsymbol{\theta}, \mathbf{x}_n) d\boldsymbol{\theta} \quad (36)$$

$$= p(\mathbf{x}_n) \quad (37)$$

so that we can apply the result 2.115 in Bishop to obtain

$$Z_n = p(\mathbf{x}_n) \quad (38)$$

$$= (1-w)\mathcal{N}(\mathbf{x}_n|\mathbf{m}^{\setminus n}, \mathbf{I} + v^{\setminus n}\mathbf{I}) + w\mathcal{N}(\mathbf{x}_n|\mathbf{0}, a\mathbf{I}) \quad (39)$$

$$= (1-w)\mathcal{N}(\mathbf{x}_n|\mathbf{m}^{\setminus n}, (1 + v^{\setminus n})\mathbf{I}) + w\mathcal{N}(\mathbf{x}_n|\mathbf{0}, a\mathbf{I}) \quad (40)$$

■