

Machine Learning 2 - Homework 5

Pascal M. Esser

May 7, 2018

Collaborators: Sindy Löwe

1

Consider a Gaussian mixture model

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

1.1

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

$$\mathbb{E}_{\text{posterior}}[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] = \sum_n^N \sum_k^K \gamma(z_{nk}) \{\ln \pi_k + \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\} \quad (2)$$

Derive update rules for $\boldsymbol{\pi}$, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$. Write the Lagrangian as:

$$\mathcal{L} = \sum_n^N \sum_k^K \gamma(z_{nk}) \{\ln \pi_k + \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\} + \lambda \left(\sum_k \pi_k - 1 \right) \quad (3)$$

Derive $\boldsymbol{\pi}_k$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \sum_n \frac{\gamma(z_{nk})}{\pi_k} + \lambda = 0 \quad (4)$$

$$\Rightarrow \pi_k = \frac{\sum_n \gamma(z_{nk})}{-\lambda} = \frac{N_k}{-\lambda} \quad (5)$$

$$\text{with } -\lambda \cdot 1 = -\lambda \sum_k \pi_k = \sum_k \sum_n \gamma(z_{nk}) = \sum_k N_k \Rightarrow \lambda = -N \quad (6)$$

$$\Rightarrow \pi_k = \frac{N_k}{N} \quad (7)$$

Derive $\boldsymbol{\mu}_k$

$$\ln \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = -\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (8)$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}_k} = \sum_n \gamma(z_{nk}) \frac{\partial \ln \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \quad (9)$$

$$= \sum_n \frac{1}{2} \gamma(z_{nk}) \frac{\partial (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)}{\partial (\mathbf{x}_n - \boldsymbol{\mu}_k)} \frac{\partial (\mathbf{x}_n - \boldsymbol{\mu}_k)}{\partial \boldsymbol{\mu}_k} \quad (10)$$

$$= \sum_n -\frac{1}{2} \gamma(z_{nk}) 2 \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) (-\mathbb{I}) \quad (11)$$

$$= \sum_n \gamma(z_{nk}) \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (12)$$

$$\boldsymbol{\Sigma}_k^{-1} \sum_n \gamma(z_{nk}) \mathbf{x}_n = \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k \sum_n \gamma(z_{nk}) \quad (13)$$

$$\Rightarrow \boldsymbol{\mu}_k = \frac{1}{N_k} \sum_n \gamma(z_{nk}) \mathbf{x}_n \quad (14)$$

Derive $\boldsymbol{\Sigma}_k$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\Sigma}_k} = \sum_n \gamma(z_{nk}) \frac{\partial \ln \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\Sigma}_k} \quad (15)$$

$$= \sum_n -\frac{1}{2} \gamma(z_{nk}) \left[\frac{\partial \ln |\boldsymbol{\Sigma}_k|}{\partial \boldsymbol{\Sigma}_k} + \frac{\partial (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)}{\partial \boldsymbol{\Sigma}_k} \right] \quad (16)$$

$$\text{using Matrix Cookbook 57 and 61} \quad (17)$$

$$= \sum_n -\frac{1}{2} \gamma(z_{nk}) \left[\boldsymbol{\Sigma}_k^{-T} - \boldsymbol{\Sigma}_k^{-\top} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-\top} \right] \quad (18)$$

$$= -\frac{1}{2} \boldsymbol{\Sigma}_k^{-\top} \sum_n \gamma(z_{nk}) - \boldsymbol{\Sigma}_k^{-\top} \sum_n \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-\top} = 0 \quad (19)$$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_n \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \quad (20)$$

1.2

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.

We don't have any changes for $\boldsymbol{\pi}_k$ and $\boldsymbol{\mu}_k$ as they do not include terms for $\boldsymbol{\Sigma}_k$.

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\Sigma}} = \sum_k \sum_n \gamma(z_{nk}) \frac{\partial \ln \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma})}{\partial \boldsymbol{\Sigma}} \quad (21)$$

$$= \sum_k \sum_n -\frac{1}{2} \gamma(z_{nk}) \left[\frac{\partial \ln |\boldsymbol{\Sigma}|}{\partial \boldsymbol{\Sigma}} + \frac{\partial (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)}{\partial \boldsymbol{\Sigma}} \right] \quad (22)$$

$$\text{using Matrix Cookbook 57 and 61} \quad (23)$$

$$= \sum_k \sum_n -\frac{1}{2} \gamma(z_{nk}) \left[\boldsymbol{\Sigma}^{-T} - \boldsymbol{\Sigma}^{-\top} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}^{-\top} \right] \quad (24)$$

$$= \sum_k -\frac{1}{2} \boldsymbol{\Sigma}^{-\top} \sum_n \gamma(z_{nk}) - \sum_k \boldsymbol{\Sigma}^{-\top} \sum_n \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}^{-\top} = 0 \quad (25)$$

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_k \sum_n \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \quad (26)$$

2

Suppose we wish to use the EM algorithm to maximize the posterior distribution $p(\Theta|\mathbf{X})$ for a model (Figure 1) containing latent variables 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

$$\sum_z p(\mathbf{Z}|\mathbf{X}, \Theta^{old}) \ln p(\mathbf{Z}, \mathbf{X}|\Theta) + \ln p(\Theta) \quad (27)$$

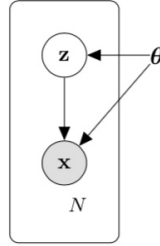


Figure 1: A simple generative model.

For adjusting the EM algorithm to maximize the posterior distribution $p(\Theta|\mathbf{X})$ we first rewrite the posterior distribution $p(\Theta|\mathbf{X})$ using Bishop [1] 9.76 and 9.77 as

$$p(\Theta|\mathbf{X}) = \frac{p(\Theta, \mathbf{X})}{p(\mathbf{X})} \quad (28)$$

$$\ln p(\Theta|\mathbf{X}) = \ln p(\Theta, \mathbf{X}) - \ln p(\mathbf{X}) \quad (29)$$

$$= \ln p(\mathbf{X}|\Theta) + \ln p(\Theta) - \ln p(\mathbf{X}) \quad (30)$$

$$= \mathcal{L}(q, \Theta) + \mathcal{D}_{KL}(q||p) + \ln p(\Theta) - \ln p(\mathbf{X}) \quad (31)$$

$$\geq \mathcal{L}(q, \Theta) + \ln p(\Theta) - \ln p(\mathbf{X}) \quad (32)$$

$$= \mathcal{L}(q, \Theta) + \ln p(\Theta) - \text{const} \quad (33)$$

$$(34)$$

This means that for an optimization with respect to q we can reuse the same E-Step as before as q only appears in $\mathcal{L}(q, \Theta)$.

Starting from Bishop [1] 9.32 and 9.33 we can furthermore rewrite the M-Step with the prior term. The original M-Step can be written as:

$$\Theta^{new} = \arg \max_{\Theta} \mathcal{Q}(\Theta, \Theta^{old}) \quad (35)$$

with

$$\mathcal{Q}(\Theta, \Theta^{old}) = \sum_z p(\mathbf{Z}|\mathbf{X}, \Theta^{old}) \ln p(\mathbf{Z}, \mathbf{X}|\Theta) \quad (36)$$

we can now rewrite this for the maximum likelihood case by adding $\ln p(\Theta)$ to the optimization as follows:

$$\mathcal{L}(q, \Theta) + \ln p(\Theta) - \text{const} = \sum_z p(\mathbf{Z}|\mathbf{X}, \Theta^{old}) \ln p(\mathbf{Z}, \mathbf{X}|\Theta) + \ln p(\Theta) - \text{const} \quad (37)$$

3

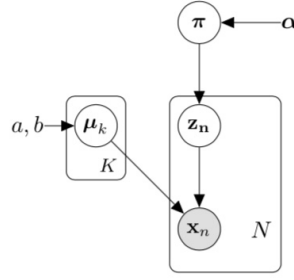


Figure 2: Mixtures of Bernoulli distribution

$$\pi | \alpha \sim \text{Dir}(\pi | \alpha) \quad (38)$$

$$z_n | \pi \sim \text{Mult}(z_n | \pi) \quad (39)$$

$$\mu_k | a_k b_k \sim \text{Beta}(\mu_k | a_k b_k) \quad (40)$$

$$\mathbf{x}_n | z_n, \mu = \{\mu_1, \dots, \mu_K\} \sim \prod_k^K (\text{Bern}(\mathbf{x}_n | \mu_k))^{z_{nk}} \quad (41)$$

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

optimize $z[\ln p(\mathbf{X}, \mathbf{Z} | \Theta)] + \ln p(\Theta)$:

$$M = \sum_n \sum_k^K \gamma(z_{nk}) \left[\ln \pi_k + \sum_i^D x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki}) \right] + \ln p(\Theta) \quad (42)$$

$$\text{with} \quad (43)$$

$$\ln p(\Theta) = \sum_k \ln \text{Beta}(\mu_k | a_k, b_k) + \ln \text{Dir}(\pi | \alpha) \quad (44)$$

$$= \sum_k \sum_i (a_k - 1) \ln \mu_{ki} + (b_k - 1) \ln(1 - \mu_{ki}) + f(a_k, b_k) + \sum_k (a_k - 1) \ln \pi_k + g(\alpha_k) \quad (45)$$

$$\frac{\partial M}{\partial \mu_{ki}} = \sum_n \gamma(z_{nk}) \left(\frac{x_{ni}}{\mu_{ki}} - \frac{1 - x_{ni}}{1 - \mu_{ki}} \right) + \frac{a_k - 1}{\mu_{ki}} - \frac{b_k - 1}{1 - \mu_{ki}} = 0 \quad (46)$$

$$\frac{1}{\mu_{ki}} \left(\sum_n \gamma(z_{nk}) x_{ni} + a_k - 1 \right) = \frac{1}{1 - \mu_{ki}} \left(\sum_n \gamma(z_{nk}) (1 - x_{ni}) + b_k - 1 \right) \quad (47)$$

$$(1 - \mu_{ki}) \left(\sum_n \gamma(z_{nk}) x_{ni} + a_k - 1 \right) = \mu_{ki} \left(\sum_n \gamma(z_{nk}) (1 - x_{ni}) + b_k - 1 \right) \quad (48)$$

$$\sum_n \gamma(z_{nk}) x_{ni} + a_k - 1 - \mu_{ki} \sum_n \gamma(z_{nk}) x_{ni} - \mu_{ki} (a_k - 1) = -\mu_{ki} \sum_n \gamma(z_{nk}) x_{ni} - \mu_{ki} \sum_n \gamma(z_{nk}) + \mu_{ki} (b_k - 1) \quad (49)$$

$$\sum_n \gamma(z_{nk}) x_{ni} + a_k - 1 = \mu_{ki} (a_k - 1 + N_k + b_k - 1) \quad (50)$$

$$\Rightarrow \mu_{ki} = \frac{\sum_n \gamma(z_{nk}) x_{ni} + a_k - 1}{a_k + N_k + b_k - 2} \quad (51)$$

alternatively we can write μ in vector form by keeping the sum over i to get $\boldsymbol{\mu}_k$ as follows (the derivation says the same with the exact of the use of the vector instead of the double index that is derived above):

$$\boldsymbol{\mu}_k = \frac{\sum_n \gamma(z_{nk}) \mathbf{x}_n + a_k - 1}{a_k + N_k + b_k - 2} \quad (52)$$

$$\frac{\partial M + \lambda (\sum_k \pi_k - 1)}{\partial \pi_k} = \sum_n \frac{\gamma(z_{nk})}{\pi_k} + \frac{\alpha_k - 1}{\pi_k} + \lambda = 0 \quad (53)$$

$$\Rightarrow \pi_k = \frac{\sum_n \gamma(z_{nk}) + \alpha_k - 1}{-\lambda} \quad (54)$$

$$= \frac{N_k + \alpha_k - 1}{-\lambda} \quad (55)$$

$$\text{with } -\lambda \sum_k \pi_k = \sum_k \sum_n \gamma(z_{nk}) + \alpha_k - K \Rightarrow -\lambda = N + \sum_k \alpha_k - K \quad (56)$$

$$\Rightarrow \pi_k = \frac{N_k + \alpha_k - 1}{N + \sum_k \alpha_k - K} \quad (57)$$

References

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006. ISBN: 0387310738.