

# 2413, Machine Learning, Tutorial 9

## Universität Bern

Simon Jenni (jenni@inf.unibe.ch)

### Expectation Maximization (EM-Algorithm)

1. Explain the differences between the Mixtures of Gaussian model (MoG) and the Gaussian Discriminant Analysis model (GDA).

**Solution.**

GDA is a supervised generative model in which we assume  $p(x|y_c) \sim \mathcal{N}(\mu_c, \Sigma_c)$ . Labels (i.e., the assignment of each training example  $x^{(i)}$  to the corresponding Gaussian) are given for the training set.

The MoG on the other hand is an unsupervised model in which we assume that each data point is sampled from a Gaussian distribution. The assignment of each  $x^{(i)}$  to one of the Gaussians is unknown in this case (i.e., we treat it as a latent variable) and has to be learnt.

2. Derive the update rule for  $\Sigma_l$  in the Maximization step (M-step) of the EM algorithm for the Mixture of Gaussian model.

**Solution.**

We need to calculate the gradient of the  $J(Q, \theta)$  with respect to  $\Sigma_l$  and set it to zero.

$$J(Q, \theta) = \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \frac{\frac{1}{(2\pi)^{\frac{1}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right)}{w_j^{(i)}}$$

Then we have:

$$\begin{aligned} \nabla_{\Sigma_l} J(Q, \theta) &= -\nabla_{\Sigma_l} \left[ \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) + \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \frac{1}{2} \log |\Sigma_j| \right] \\ &= -\nabla_{\Sigma_l} \frac{1}{2} \sum_{i=1}^m w_l^{(i)} \left[ (x^{(i)} - \mu_l)^T \Sigma_l^{-1} (x^{(i)} - \mu_l) + \log |\Sigma_l| \right] \end{aligned}$$

We set  $\Lambda_l = \Sigma_l^{-1}$  and solve for  $\Lambda_l$ .

$$= \nabla_{\Lambda_l} \left[ -\frac{1}{2} \sum_{i=1}^m w_l^{(i)} (x^{(i)} - \mu_l)^T \Lambda_l (x^{(i)} - \mu_l) + \frac{1}{2} \sum_{i=1}^m w_l^{(i)} \log |\Lambda_l| \right] \quad (1)$$

$$= \nabla_{\Lambda_l} \left[ -\frac{1}{2} \sum_{i=1}^m w_l^{(i)} \text{tr}((x^{(i)} - \mu_l)^T \Lambda_l (x^{(i)} - \mu_l)) + \frac{1}{2} \sum_{i=1}^m w_l^{(i)} \log |\Lambda_l| \right] \quad (2)$$

$$= \nabla_{\Lambda_l} \left[ -\frac{1}{2} \sum_{i=1}^m w_l^{(i)} \text{tr}((x^{(i)} - \mu_l)(x^{(i)} - \mu_l)^T \Lambda_l) + \frac{1}{2} \sum_{i=1}^m w_l^{(i)} \log |\Lambda_l| \right] \quad (3)$$

$$= -\frac{1}{2} \sum_{i=1}^m w_l^{(i)} (x^{(i)} - \mu_l)(x^{(i)} - \mu_l)^T + \frac{1}{2} \Lambda_l^{-T} \sum_{i=1}^m w_l^{(i)} \quad (4)$$

Where

- (2) follows from (1) by  $a = \text{tr}(a)$ ,  $\forall a \in \mathcal{R}$ . Note that  $(x^{(i)} - \mu_j)^T \Lambda_l (x^{(i)} - \mu_j) \in \mathcal{R}$ .
- (3) follows from (2) by  $\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$ .
- (4) follows from (3) by  $\nabla_A \text{tr}(BA) = B^T$  and  $\nabla_A \log |A| = A^{-T}$ . Note that  $((x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T)^T = (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T$ .

By setting to zero we have:

$$\Lambda_l^{-T} = \Lambda_l^{-1} = \Sigma_l = \frac{\sum_{i=1}^m w_l^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_l^{(i)}}$$

## Factor Analysis

3. Assume that  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$  are sampled i.i.d. from a distribution described by the factor analysis model

$$z \sim \mathcal{N}(0, I) \quad (5)$$

$$\epsilon \sim \mathcal{N}(0, \Psi) \quad (6)$$

$$x = \mu + \Lambda z + \epsilon. \quad (7)$$

What is the optimal  $\mu$ ? Use Maximum-Likelihood estimation.

**Solution.**

The samples are drawn from the distribution  $x \sim \mathcal{N}(\mu, \Lambda\Lambda^T + \Psi)$ . The log-likelihood function according to the ML estimate is

$$l(\mu) = \log \prod_{i=1}^m \frac{\exp(-\frac{1}{2}(x^{(i)} - \mu)^T(\Lambda\Lambda^T + \Psi)^{-1}(x^{(i)} - \mu))}{(2\pi)^{n/2}|\Lambda\Lambda^T + \Psi|^{1/2}} \quad (8)$$

$$= \sum_{i=1}^m -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\Lambda\Lambda^T + \Psi|) + \quad (9)$$

$$\sum_{i=1}^m -\frac{1}{2} (x^{(i)} - \mu)^T (\Lambda\Lambda^T + \Psi)^{-1} (x^{(i)} - \mu) \quad (10)$$

Note that the negative log-likelihood is a convex quadratic function in  $\mu$ , therefore we can find the optimal  $\mu$  if we set the gradient to 0. The gradient of the log-likelihood w.r.t.  $\mu$  is

$$\nabla_{\mu} l(\mu) = \nabla_{\mu} \sum_{i=1}^m -\frac{1}{2} (x^{(i)} - \mu)^T (\Lambda\Lambda^T + \Psi)^{-1} (x^{(i)} - \mu) \quad (11)$$

$$= \sum_{i=1}^m -(\Lambda\Lambda^T + \Psi)^{-1} \mu + (\Lambda\Lambda^T + \Psi)^{-1} x^{(i)}. \quad (12)$$

From here, the solution is not very surprisingly,

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}. \quad (13)$$