

2413, Machine Learning, Tutorial 10

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 Σ_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 Σ_l and set it to zero.

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

It will be easier to maximize with respect to $\Lambda_l = \Sigma_l^{-1}$ instead of Σ_l . Let us therefore calculate:

$$\nabla_{\Lambda_l} J(Q, \theta) = \nabla_{\Lambda_l} \left[\sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \left(-\frac{1}{2} \log |\Lambda_j^{-1}| - \frac{1}{2} (x^{(i)} - \mu_j)^T \Lambda_j (x^{(i)} - \mu_j) \right) \right] \quad (1)$$

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

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

$$= \frac{1}{2} \sum_{i=1}^m w_l^{(i)} \nabla_{\Lambda_l} \left[\log |\Lambda_l| - \text{tr}((x^{(i)} - \mu_l)(x^{(i)} - \mu_l)^T \Lambda_l) \right] \quad (4)$$

$$= \frac{1}{2} \sum_{i=1}^m w_l^{(i)} \left(\Lambda_l^{-T} - (x^{(i)} - \mu_l)(x^{(i)} - \mu_l)^T \right) \quad (5)$$

Where

- (1) follows by substituting Λ_j and ignoring all terms not containing Λ_j .
- (2) follows from (1) by ignoring all terms not containing Λ_l and using $|A^{-1}| = 1/|A|$.
- (3) follows from (2) by $a = \text{tr}(a)$, $\forall a \in \mathbb{R}$. Note that $(x^{(i)} - \mu_j)^T \Lambda_l (x^{(i)} - \mu_j) \in \mathbb{R}$.
- (4) follows from (3) by $\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$.
- (5) follows from (4) 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

$$z \sim \mathcal{N}(0, I) \quad (6)$$

$$\epsilon \sim \mathcal{N}(0, \Psi) \quad (7)$$

$$x = \mu + \Lambda z + \epsilon. \quad (8)$$

What is the optimal μ ? 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

$$l(\mu) = \log \prod_{i=1}^m \frac{\exp(-\frac{1}{2}(x^{(i)} - \mu)^T(\Lambda\Lambda^T + \Psi)^{-1}(x^{(i)} - \mu))}{(2\pi)^n |\Lambda\Lambda^T + \Psi|^{1/2}} \quad (9)$$

$$= \sum_{i=1}^m -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\Lambda\Lambda^T + \Psi|) \quad (10)$$

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

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

$$\nabla_{\mu} l(\mu) = \nabla_{\mu} \sum_{i=1}^m -\frac{1}{2} (x^{(i)} - \mu)^T (\Lambda\Lambda^T + \Psi)^{-1} (x^{(i)} - \mu) \quad (12)$$

$$= \sum_{i=1}^m -(\Lambda\Lambda^T + \Psi)^{-1} \mu + (\Lambda\Lambda^T + \Psi)^{-1} x^{(i)}. \quad (13)$$

From here, the solution is not very surprisingly,

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}. \quad (14)$$