

EM algorithm for Gaussian mixtures

Kitov Victor

1 Lemmas

Lemma 1. For any square non-degenerate matrix B $\frac{\partial |B|}{\partial B} = |B|B^{-1}$

Proof. Row expansion of determinant:

$$|B| = \sum_{j=1}^D b_{ij} A_{ij}$$

where A_{ij} is algebraic complement (also known as cofactor) of B for b_{ij} . So

$$\frac{\partial |B|}{\partial b_{ij}} = A_{ij}$$

Combining this with the fact that $B^{-1} = \frac{1}{|B|}A^T$, where $\{A\}_{ij} = A_{ij}$, obtain that $\frac{\partial |B|}{\partial B} = |B|B^{-1}$. □

Lemma 2. For any square symmetric non-degenerate matrix B $\frac{\partial \ln |B|}{\partial B} = B^{-1}$

Proof. $\frac{\partial \ln |B|}{\partial B} = \frac{1}{|B|} |B| B^{-1} = B^{-1}$ since B is symmetric. □

Lemma 3. For any square matrix U $\frac{\partial \{a^T U b\}}{\partial U} = ab^T$.

Proof.

$$\begin{aligned} a^T U b &= \sum_{i,j} a_i u_{ij} b_j \\ \frac{\partial \{a^T U b\}}{\partial u_{ij}} &= a_i b_j \\ \frac{\partial \{a^T U b\}}{\partial U} &= ab^T \end{aligned}$$

□

2 Definitions

Training set x_1, x_2, \dots, x_N .

If component is known

$$p(x|z, \theta) = N(x|\mu_z, \Sigma_z) = \frac{1}{(2\pi)^{D/2} |\Sigma_z|^{1/2}} e^{-\frac{1}{2}(x-\mu_z)^T \Sigma_z^{-1} (x-\mu_z)} \quad (1)$$

Priors

$$p(z|\theta) = \phi_z \quad (2)$$

Each object is modelled with mixture of Z Gaussians with parameters $\{\mu_z, \Sigma_z\}$ and priors ϕ_z , $z = 1, 2, \dots, Z$.

$$p(x|\theta) = \sum_{z=1}^Z \phi_z N(x|\mu_z, \Sigma_z) = \sum_{z=1}^Z \phi_z \frac{1}{(2\pi)^{D/2} |\Sigma_z|^{1/2}} e^{-\frac{1}{2}(x-\mu_z)^T \Sigma_z^{-1} (x-\mu_z)} \quad (3)$$

Parameters of mixture: $\theta = \{\phi_z, \mu_z, \Sigma_z\}_{z=1}^Z$.

3 Expectation step

$$q_n(z) = w_{nz} = P(z|x_n) = \frac{P(z, x_n)}{P(x_n)} = \frac{P(z, x_n)}{\sum_k P(k, x_n)} = \frac{P(z)P(x_n|z)}{\sum_k P(k)P(x_n|k)} = \frac{\hat{\phi}_z N(x_n|\hat{\mu}, \hat{\Sigma}_z)}{\sum_k \hat{\phi}_k N(x_n|\hat{\mu}_k, \hat{\Sigma}_k)}$$

4 Maximization step

4.1 Optimization task

Given, that $\phi_z = p(z)$, be need to ensure

$$\sum_{z=1}^Z \phi_z = 1 \tag{4}$$

Given (4) we need to solve

$$\sum_{n=1}^N \sum_{z=1}^Z w_{nz} \ln \frac{p(x_n, z|\theta)}{w_{nz}} \rightarrow \max_{\theta}$$

This is equivalent to

$$\sum_{n=1}^N \sum_{z=1}^Z w_{nz} \ln p(x_n, z|\theta) = \sum_{n=1}^N \sum_{z=1}^Z w_{nz} \ln p(z|\theta)p(x_n|z, \theta) \rightarrow \max_{\theta} \tag{5}$$

Substituting (2) and (1) into (5) and using (4) we obtain final optimization task for the M-step:

$$\begin{cases} \sum_{n=1}^N \sum_{z=1}^Z w_{nz} \ln \phi_z + \sum_{n=1}^N \sum_{z=1}^Z w_{nz} \ln (2\pi)^{-D/2} + \frac{1}{2} \sum_{n=1}^N \sum_{z=1}^Z w_{nz} \ln |\Sigma_z^{-1}| \\ - \frac{1}{2} \sum_{n=1}^N \sum_{z=1}^Z w_{nz} (x_n - \mu_z)^T \Sigma_z^{-1} (x_n - \mu_z) \rightarrow \max_{\theta} \\ \sum_{z=1}^Z \phi_z = 1 \end{cases} \tag{6}$$

4.2 Optimization for ϕ_z .

Optimization for ϕ_z leads to Lagrangian

$$\mathcal{L} = \sum_{n=1}^N \sum_{z=1}^Z w_{nz} \ln \phi_z + \lambda (1 - \sum_{z=1}^Z \phi_z) \rightarrow \text{extr}$$

Optimal values satisfy

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \phi_z} &= \sum_{n=1}^N w_{nz} \frac{1}{\phi_z} - \lambda = 0 \\ \phi_z &= \frac{1}{\lambda} \sum_{n=1}^N w_{nz} \end{aligned}$$

Using (4) we find λ :

$$\begin{aligned} 1 &= \sum_z \phi_z = \frac{1}{\lambda} \sum_n \sum_z w_{nz} \\ \lambda &= \sum_n \sum_z w_{nz} = \sum_{n=1}^N 1 = N \end{aligned}$$

It follows that

$$\phi_z = \frac{1}{N} \sum_{n=1}^N w_{nz}$$

Interpretation: ϕ_z is sample average of component correspondences, when each object may partly belong to each component.

4.3 Optimization for μ_z .

Optimization for μ_z leads to

$$\begin{aligned}\mathcal{L} &= \sum_{n=1}^N \sum_{z=1}^Z w_{nz} (x_n - \mu_z)^T \Sigma_z^{-1} (x_n - \mu_z) \rightarrow \max_{\mu_z} \\ \frac{\partial \mathcal{L}}{\partial \mu_z} &= 2 \sum_n w_{nz} \Sigma_z^{-1} (x_n - \mu_z) (-1) = 0 \\ \sum_n w_{nz} \Sigma_z^{-1} (x_n - \mu_z) &= \Sigma_z^{-1} \sum_n \{w_{nz} (x_n - \mu_z)\} = 0 \\ \sum_n \{w_{nz} (x_n - \mu_z)\} &= \sum_n w_{nz} x_n - \sum_n w_{nz} \mu_z = 0 \\ \mu_z &= \frac{\sum_n w_{nz} x_n}{\sum_n w_{nz}}\end{aligned}$$

Interpretation: μ_z equals to weighted average of objects with weights=correspondences to corresponding component.

4.4 Optimization for Σ

From (6) we obtain the following optimization problem to find covariance matrices:

$$\sum_{n=1}^N \sum_{z=1}^Z w_{nz} \ln |\Sigma_z^{-1}| - \sum_{n=1}^N \sum_{z=1}^Z w_{nz} (x_n - \mu_z)^T \Sigma_z^{-1} (x_n - \mu_z) \rightarrow \max_{\Sigma_1, \dots, \Sigma_Z} \quad (7)$$

4.4.1 Unconstrained Σ_z .

For each individual Σ_z we get

$$\mathcal{L} = \sum_{n=1}^N w_{nz} \ln |\Sigma_z^{-1}| - \sum_{n=1}^N w_{nz} (x_n - \mu_z)^T \Sigma_z^{-1} (x_n - \mu_z) \rightarrow \max_{\Sigma_z} \quad (8)$$

Using lemmas (3), (2)

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \Sigma_z^{-1}} &= \sum_{n=1}^N \{w_{nz} \Sigma_z\} - \sum_{n=1}^N w_{nz} (x_n - \mu_z)(x_n - \mu_z)^T = 0 \\ \Sigma_z &= \frac{\sum_{n=1}^N w_{nz} (x_n - \mu_z)(x_n - \mu_z)^T}{\sum_{n=1}^N w_{nz}}\end{aligned}$$

Interpretation: Σ_z equals to weighted sample covariance matrix with weights=correspondences of objects to corresponding component.

4.4.2 Constraint: common covariance $\Sigma_z = \Sigma \forall z$.

4.4.3 Constraint: diagonal covariance $\Sigma_z = \text{diag}\{\sigma_{z,1}^2, \dots, \sigma_{z,D}^2\} \forall z$.

4.4.4 Constraint: spherical covariance $\Sigma_z = \sigma_z^2 I \forall z$.