## Machine Learning for Computer Vision

January 23, 2019 Topic: Clustering

## Exercise 1: K-Means Compression (Programming)

a) K-Means finds the parameters (cluster means) that minimize the assignment cost:

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{N} \sum_{j=1}^{K} A_{ij} ||\boldsymbol{x}_i - \boldsymbol{\mu}_j||^2 , \qquad (1)$$

where  $\boldsymbol{\theta} = \{\boldsymbol{\mu}_k : 1 \leq k \leq K\}$  and  $A_{ij}$  is 1 if  $x_i$  is assigned to cluster j and 0 otherwise.

b) See code.

## Exercise 2: Expectation-Maximization for GMM (Programming)

a) EM finds the parameters (means, covariances and mixture coefficients) that maximize the conditional data log-likelihood:

$$\boldsymbol{\theta}^* = \arg\max_{\boldsymbol{\theta}} \log p(X|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})$$
(2)

where  $\boldsymbol{\theta} = \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \boldsymbol{\pi}_k : 1 \leq k \leq K\}$ 

b) See code.

## Exercise 3: Expectation-Maximization for GMM

In the standard EM algorithm, we first define the responsibilities  $\gamma$  as

$$\gamma_{nk} = p(z_{nk} = 1 | x_n) = \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)} \quad , z_{nk} \in \{0, 1\}, \sum_{k=1}^K z_{nk} = 1$$

a) Find the optimal means, covariances and mixing coefficients that maximize the data likelihood. How can we interpret the results?

We want to maximize the data likelihood, so as usual we minimize the negative log-likelihood:

$$-\mathcal{L}\mathcal{L} = -\log p(X|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = -\log \prod_{n} \sum_{k} \pi_{k} \mathcal{N}(x_{n}|\boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})$$
(3)

This time we minimize 3 times independently with respect to the means, the covariances and the mixture coefficients:

$$\mu_k^* = \underset{\mu_k}{\operatorname{arg\,min}} -\mathcal{L}\mathcal{L} \tag{4}$$

$$\Sigma_k^* = \underset{\Sigma_k}{\arg\min} - \mathcal{LL}$$
(5)

$$\pi_k^* = \underset{\pi_k}{\operatorname{arg\,min}} -\mathcal{L}\mathcal{L} \tag{6}$$

In the following, to avoid confusion of sums and covariances, we denote covariance  $\Sigma_k$  as  $C_k$ . To simplify some expressions, let us agree on the following notation:

$$\mathcal{N}_{nk} \equiv \mathcal{N}(x_n | \mu_k, C_k) \tag{7}$$

$$Z_k \equiv ((2\pi)^d |C_k|)^{1/2}$$
(8)

$$D_{nk} \equiv (x_n - \mu_k)^T C_k^{-1} (x_n - \mu_k)$$
(9)

Therefore 
$$\mathcal{N}_{nk} = Z_k^{-1} \exp\{-\frac{1}{2}D_{nk}\}$$
 (10)

Thus, we have:

$$-\mathcal{L}\mathcal{L} = -\sum_{n} \log \sum_{k} \pi_{k} \mathcal{N}_{nk}$$
$$= -\sum_{n} \log \sum_{k} \pi_{k} Z_{k}^{-1} \exp(-\frac{1}{2} D_{nk})$$

Solving for the means:

$$\frac{\partial \mathcal{LL}}{\partial \mu_k} = \sum_n \frac{1}{\sum_j \pi_j \mathcal{N}_{nj}} \frac{\partial \sum_k \pi_k Z_k^{-1} \exp(-\frac{1}{2} D_{nk})}{\partial \mu_k}$$
(11)

$$=\sum_{n}\frac{1}{\sum_{j}\pi_{j}\mathcal{N}_{nj}}\pi_{k}Z_{k}^{-1}\frac{\partial\exp(-\frac{1}{2}D_{nk})}{\partial\mu_{k}}$$
(12)

$$=\sum_{n}\frac{1}{\sum_{j}\pi_{j}\mathcal{N}_{nj}}\pi_{k}Z_{k}^{-1}\exp(-\frac{1}{2}D_{nk})C_{k}^{-1}(x_{n}-\mu_{k})$$
(13)

$$=\sum_{n} \frac{\pi_k \mathcal{N}_{nk}}{\sum_j \pi_j \mathcal{N}_{nj}} C_k^{-1} (x_n - \mu_k)$$
(14)

$$=\sum_{n}\gamma_{nk}C_{k}^{-1}(x_{n}-\mu_{k})$$
(15)

(16)

Setting  $-\frac{\partial \mathcal{LL}}{\partial \mu_k} \stackrel{!}{=} 0$  gives us:

$$\sum_{n} \gamma_{nk} C_k^{-1} \mu_k = \sum_{n} \gamma_{nk} C_k^{-1} x_n \tag{17}$$

$$C_k^{-1}\mu_k \sum_n \gamma_{nk} = C_k^{-1} \sum_n \gamma_{nk} x_n \tag{18}$$

$$C_k^{-1}\mu_k \sum_n \gamma_{nk} = C_k^{-1} \sum_n \gamma_{nk} x_n \tag{19}$$

$$\mu_k \sum_n \gamma_{nk} = \sum_n \gamma_{nk} x_n \tag{20}$$

$$\mu_k = \frac{\sum_n \gamma_{nk} x_n}{\sum_n \gamma_{nk}} \tag{21}$$

Solving for the covariances:

$$\frac{\partial \mathcal{LL}}{\partial C_k} = \sum_n \frac{1}{\sum_j \pi_j \mathcal{N}_{nj}} \frac{\partial \sum_k \pi_k Z_k^{-1} \exp(-\frac{1}{2}D_{nk})}{\partial C_k}$$
(22)

$$=\sum_{n}\frac{1}{\sum_{j}\pi_{j}\mathcal{N}_{nj}}\pi_{k}\frac{\partial Z_{k}^{-1}\exp(-\frac{1}{2}D_{nk})}{\partial C_{k}}$$
(23)

$$=\sum_{n}\frac{1}{\sum_{j}\pi_{j}\mathcal{N}_{nj}}\pi_{k}\left(\frac{\partial Z_{k}^{-1}}{\partial C_{k}}\exp(-\frac{1}{2}D_{nk})+Z_{k}^{-1}\frac{\partial\exp(-\frac{1}{2}D_{nk})}{\partial C_{k}}\right)$$
(24)

$$=\sum_{n} \frac{1}{\sum_{j} \pi_{j} \mathcal{N}_{nj}} \pi_{k} \left( \left( -\frac{1}{2} Z_{k}^{-1} C_{k}^{-1} \right) \exp\left(-\frac{1}{2} D_{nk}\right) + \frac{1}{2} Z_{k}^{-1} \exp\left(-\frac{1}{2} D_{nk}\right) C_{k}^{-1} (x_{n} - \mu_{k}) (x_{n} - \mu_{k})^{T} C_{k}^{-1} \right)$$

$$(25)$$

$$= (-\frac{1}{2}) \sum_{n} \frac{1}{\sum_{j} \pi_{j} \mathcal{N}_{nj}} \pi_{k} Z_{k}^{-1} \exp(-\frac{1}{2} D_{nk}) \left( C_{k}^{-1} - C_{k}^{-1} (x_{n} - \mu_{k}) (x_{n} - \mu_{k})^{T} C_{k}^{-1} \right)$$
(26)

$$= \left(-\frac{1}{2}\right) \sum_{n} \gamma_{nk} \left(C_k^{-1} - C_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^T C_k^{-1}\right)$$
(27)

(28)

Here, we used the derivative of the determinant as follows:

$$\frac{\partial Z_k^{-1}}{\partial C_k} = \frac{\partial ((2\pi)^d |C_k|)^{-\frac{1}{2}}}{\partial C_k} = ((2\pi)^d)^{-\frac{1}{2}} \frac{\partial (|C_k|)^{-\frac{1}{2}}}{\partial C_k}$$
(29)

$$= ((2\pi)^d)^{-\frac{1}{2}} (-\frac{1}{2}) |C_k|^{-\frac{3}{2}} \frac{\partial (|C_k|)}{\partial C_k} = ((2\pi)^d)^{-\frac{1}{2}} (-\frac{1}{2}) |C_k|^{-\frac{3}{2}} |C_k| (C_k^{-1})^T$$
(30)

$$= (-\frac{1}{2})((2\pi)^d)^{-\frac{1}{2}}|C_k|^{-\frac{1}{2}}C_k^{-1} = -\frac{1}{2}Z_k^{-1}C_k^{-1}$$
(31)

and the derivative of the Mahalanobis distance as:

$$\frac{\partial x^T C^{-1} x}{\partial C} = -C^{-T} x x^T C^{-T} = -C^{-1} x x^T C^{-1}$$
(32)

Setting  $-\frac{\partial \mathcal{LL}}{\partial C_k} \stackrel{!}{=} 0$  gives us:

$$\sum_{n} \gamma_{nk} C_k^{-1} = \sum_{n} \gamma_{nk} C_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^T C_k^{-1}$$
(33)

$$C_k^{-1} \sum_n \gamma_{nk} = C_k^{-1} \sum_n \gamma_{nk} (x_n - \mu_k) (x_n - \mu_k)^T C_k^{-1}$$
(34)

$$\sum_{n} \gamma_{nk} = \sum_{n} \gamma_{nk} (x_n - \mu_k) (x_n - \mu_k)^T C_k^{-1}$$
(35)

$$C_k = \frac{\sum_n \gamma_{nk} (x_n - \mu_k) (x_n - \mu_k)^T}{\sum_n \gamma_{nk}}$$
(36)

Solving for the mixture coefficients: Here we must take into account that  $\sum_k \pi_k = 1$ . We enforce this constraint with a Lagrange multiplier. Our objective then becomes:

$$\mathcal{LL}' = \mathcal{LL} + \lambda(\sum_{k} \pi_k - 1)$$
(37)

where  $\lambda < 0$ .

Deriving w.r.t.  $\pi_k$ , we get

$$\frac{\partial \mathcal{LL}'}{\partial \pi_k} = \sum_n \frac{1}{\sum_j \pi_j \mathcal{N}_{nj}} \frac{\partial \sum_k \pi_k \mathcal{N}_{nk}}{\partial \pi_k} + \lambda$$
(38)

$$=\sum_{n}\frac{1}{\sum_{j}\pi_{j}\mathcal{N}_{nj}}\mathcal{N}_{nk}+\lambda$$
(39)

$$=\sum_{n}\frac{\gamma_{nk}}{\pi_{k}}+\lambda\tag{40}$$

Setting equal to zero and solving for  $\lambda$ , we get

$$\lambda = -\sum_{n} \frac{\gamma_{nk}}{\pi_k} \tag{41}$$

$$\lambda \pi_k = -\sum_n \gamma_{nk} \tag{42}$$

$$\sum_{k} \lambda \pi_k = -\sum_{k} \sum_{n} \gamma_{nk} \tag{43}$$

$$\lambda = -N \tag{44}$$

Now we can plug this back to the objective and actually solve for  $\pi_k$ :

$$\frac{\partial \mathcal{L}\mathcal{L}'}{\partial \pi_k} = \sum_n \frac{\gamma_{nk}}{\pi_k} - N \stackrel{!}{=} 0 \tag{45}$$

$$\frac{1}{\pi_k} \sum_n \gamma_{nk} = N \tag{46}$$

$$\pi_k = \frac{\sum_n \gamma_{nk}}{N} = \frac{N_k}{N} \tag{47}$$

We can interpret these results as weighted averages of means and covariances, the weights corresponding to the responsibilities  $\gamma_{nk}$ . The mixture coefficients  $\pi_k$  are simply the ratio of data points explained by each component.

b) Define the complete-data-log-likelihood. What is the difference to the standard log-likelihood?

Assuming we observe not only the data but also the binary latent variables Z we define the complete data likelihood as:

$$p(X, Z | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{C}) = \prod_{n} p(z_n | \boldsymbol{\pi}) p(x_n | z_n, \boldsymbol{\mu}, \boldsymbol{C})$$
(48)

where  $p(z_n|\boldsymbol{\pi}) = \prod_k \pi_k^{z_{nk}}$  and  $p(x_n|z_n, \boldsymbol{\mu}, \boldsymbol{C}) = \prod_k \mathcal{N}(x_n|\mu_k, C_k)^{z_{nk}}$ . Remember that  $\sum_k z_{nk} = 1$ .

Since now we only have products, we can more easily compute the logarithm:

$$\log p(X, Z | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{C}) = \sum_{n} \sum_{k} z_{nk} (\log \pi_k + \log \mathcal{N}(x_n | \boldsymbol{\mu}_k, C_k))$$
(49)

Of course in practice, the latent variables are not known, so we maximize the *expectation*:

$$\mathbb{E}[\log p(X, Z | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{C})] = \sum_{n} \sum_{k} \mathbb{E}[z_{nk}](\log \pi_k + \log \mathcal{N}(x_n | \boldsymbol{\mu}_k, C_k))$$
(50)

where we know that  $\mathbb{E}[z_{nk}] = \gamma_{nk}$ .

The theory says that the log-marginal is also maximized implicitly!