

PRML chapter9

1. K-means clustering

- our purpose is to deivide $\{x_1, x_2, \dots, x_N\}$ into K groups
- Objective function is

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

- $\boldsymbol{\mu}_k$ is a center of the cluster, and r_{nk} is what cluster \mathbf{x}_n in
- First we initialize $\boldsymbol{\mu}_k$, then minimize J with respect to the r_{nk} , keeping the $\boldsymbol{\mu}_k$ fixed. And then minimize J with respect to the $\boldsymbol{\mu}_k$, keeping r_{nk} fixed.
- J can be written ad follows

$$J = \sum_{k=1}^K r_{1k} \|x_1 - \mu_k\|^2 + \dots + \sum_{k=1}^K r_{Nk} \|x_N - \mu_k\|^2$$
$$r_{nk} = \begin{cases} 1 & k = \arg \min_j \|x_n - \mu_j\|^2 \\ 0 & \end{cases}$$

- minimize J with respect to the $\boldsymbol{\mu}_k$, keeping r_{nk} fixed.

$$\frac{\partial J}{\partial \boldsymbol{\mu}_k} = \frac{\partial}{\partial \boldsymbol{\mu}_k} \left(\sum_{n=1}^N r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 \right) = 2 \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0$$

- gain

$$\boldsymbol{\mu}_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}$$

- K-medoids algorithm is to make objective function

$$\tilde{J} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \mathcal{V}(\mathbf{x}_n, \boldsymbol{\mu}_k)$$

2. Mixtures of Gaussians

- the conditional distribution of \mathbf{x} given a particular value for \mathbf{z} is

$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$$

-

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

- we are able to work with the joint distribution $p(\mathbf{x}, \mathbf{z})$ instead of the marginal distribution $p(\mathbf{x})$ which is written in chapter 2.9

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} \prod_{k=1}^K N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} = \prod_{k=1}^K (\pi_k N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{z_k}$$

- the conditional probability of \mathbf{z} given \mathbf{x} plays an important role

$$\gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) = \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

(a) Maximum likelihood

- the log of the likelihood function is

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

- there is a problem that distribution might be too shape around the mean to cause overfitting \rightarrow Bayesian approach prevent it

(b) EM for Gaussian mixtures

-

$$0 = \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{\partial}{\partial \boldsymbol{\mu}_k} \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) = \sum_{n=1}^N \frac{\pi_k N(\mathbf{x}_n|\boldsymbol{\mu}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

- then we gain

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

(under $N_k = \sum_{n=1}^N \gamma(z_{nk})$)

- derivate it with respect to $\boldsymbol{\Sigma}_k$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

- maximizw it with respect to the mixing coefficient π_k , using lagrangian function

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

we gain

$$\pi_k = \frac{N_k}{N}$$

- repeat this three steps until convergence criterion is satisfied (EM algorithm)

3. An Alternative View of EM

- first, we chose an initial value of parameters $\boldsymbol{\theta}^{old}$
- Second, we take, **Estep** $Evaluate p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old})$
- let Q be

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$$

and take M step,

$$\boldsymbol{\theta}^{new} = \arg \max \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{old})$$

- repeat this step until convergence criterion is satisfied

(a) Gaussian mixtures revisited

- the likelihood for the complete data set \mathbf{X}, \mathbf{Z}

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_{nk}}$$

- posterior distribution of \mathbf{Z} is

$$p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) \propto \prod_{n=1}^N \prod_{k=1}^K [\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_{nk}}$$

- we gain

$$E[z_{nk}] = \frac{\sum_{z_{nk}} z_{nk} [\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_{nk}}}{\sum_{z_{nj}} [\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)]^{z_{nj}}} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} = \gamma(z_{nk})$$

(b) Relation to K-means

- EM algorithm makes a soft assignment based on the posterior probabilities (K-means are hard)

(c) Mixtures of Bernoulli distributions

- Consider the a finite mixture of Bernoulli distributions. let \mathbf{x} be a variable, $\boldsymbol{\mu}$ be a parameter, then

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

- the log likelihood function is

$$\ln p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\pi}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k) \right\}$$

- prior distribution for the latent value is

$$p(\mathbf{z} | \boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{z_k}$$

- since

$$p(x | z, \boldsymbol{\mu}) = \prod_{k=1}^K p(x | \boldsymbol{\mu}_k)^{z_k} = \prod_{k=1}^K \left(\prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i} \right)^{z_k}$$

, we could gain

$$p(\mathbf{x}, \mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\pi}) = p(\mathbf{x} | \mathbf{z}, \boldsymbol{\mu}) p(\mathbf{z} | \boldsymbol{\pi}) = \prod_{k=1}^K (\pi_k p(\mathbf{x} | \boldsymbol{\mu}_k))^{z_k} = \prod_{k=1}^K \left(\pi_k \left(\prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i} \right) \right)^{z_k}$$

- using

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\pi}) = \sum_{n=1}^N \ln p(\mathbf{x}_n, \mathbf{z}_n | \boldsymbol{\mu}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \ln \pi_k \left(\prod_{t=1}^D \mu_{kt}^{x_{nt}} (1 - \mu_{kt})^{1-x_{nt}} \right),$$

$$E_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{ \ln \pi_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln (1 - \mu_{ki})] \}$$

- derivate it with respect to $\boldsymbol{\mu}_k$ and π_k

4. The EM algorithm in general

- Here we give a very general treatment of the EM algorithm and in the process provide a proof that the EM algorithm derived heuristically in Sections 9.2 and 9.3 for Gaussian mixtures does indeed maximize the likelihood function
- consider

$$p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$$

- we could decompose log likelihood as

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q\|p),$$

under

$$\text{KL}(q\|p) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}$$

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right\}$$