

# Homework 5

1(a) Given

$$l(\theta) = \sum_{i=1}^n \log \left[ \sum_{j=1}^K p(x^{(i)} | \mu_j, \sigma_j^2) p(j) \right]$$

$$\frac{\partial}{\partial \mu_j} l(\theta) = \sum_{i=1}^n \frac{p(j) N(x^{(i)} | \mu_j, \sigma_j^2) \nabla_{\mu_j} \left\{ -\frac{1}{2} (x^{(i)} - \mu_j)^T \sigma_j^{-1} (x^{(i)} - \mu_j) \right\}}{\sum_{j'=1}^K p(j') N(x^{(i)} | \mu_{j'}, \sigma_{j'}^2)}$$

$$= \sum_{i=1}^n p_{i,j} \sigma_j^{-1} (x^{(i)} - \mu_j)$$

$$1(b) \frac{\partial}{\partial p(j)} l(\theta) = \sum_{i=1}^n \frac{N(x^{(i)} | \mu_j, \sigma_j^2)}{\sum_{j'=1}^K p(j') N(x^{(i)} | \mu_{j'}, \sigma_{j'}^2)}$$

$$= \frac{1}{p(j)} \sum_{i=1}^n p_{i,j} \text{ [Given } p_{i,j} = p(j | x^{(i)})]$$

No, it will not be a valid probability distribution

$$1(c) \frac{\partial}{\partial P(j)} \mathcal{L}(\theta) = \frac{1}{P(j)} \sum_{i=1}^n P_{i,j}$$

$$P(j) = \frac{\exp(\omega_j^0)}{\sum_{j'=1}^K \exp(\omega_{j'}^0)}$$

using Lagrange multiplier,

$$P(j) = \frac{\sum_{i=1}^n P_{i,j}}{\lambda}$$

Summing over  $j$  and normalizing,

$$P(j) = \frac{\sum_{i=1}^n P_{i,j}}{N}$$

$$\Rightarrow \frac{\sum_{i=1}^n P_{i,j}}{P(j)} = N$$

$$\therefore \frac{\partial}{\partial \omega_j} \mathcal{L}(\theta) \propto \sum_{i=1}^n P_{i,j} - P(j)$$

Problem 1(d):

$$\frac{\partial}{\partial \mathbf{c}_j} \ell(\theta) = \frac{n}{\sum_{i=1}^n} \frac{p(i) \nabla \mathbf{c}_j N(\mathbf{x}^{(i)} | \mu_j, \mathbf{c}_j)}{p(i') N(\mathbf{x}^{(i')} | \mu_{j'}, \mathbf{c}_{j'})}$$

$$\begin{aligned} \text{Here } \nabla \mathbf{c}_j N(\mathbf{x} | \mu_j, \mathbf{c}_j) &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\mathbf{c}_j|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu_j)^T \mathbf{c}_j^{-1} (\mathbf{x} - \mu_j)\right\} \\ &\quad \nabla \mathbf{c}_j \cdot \left\{ \nabla \mathbf{c}_j \left(-\frac{1}{2}(\mathbf{x} - \mu_j)^T \mathbf{c}_j^{-1} (\mathbf{x} - \mu_j)\right) - \right. \\ &\quad \left. \mathbf{c}_j^{-1} \nabla \mathbf{c}_j |\mathbf{c}_j| \right\} \\ &= N(\mathbf{x} | \mu_j, \mathbf{c}_j) \nabla (\log N(\mathbf{x} | \mu_j, \mathbf{c}_j)) \end{aligned}$$

$$\text{So, } \mathbf{c}_j = \frac{\sum_{i=1}^n p_{i,j} (\mathbf{x}^{(i)} - \mu_j)(\mathbf{x}^{(i)} - \mu_j)^T}{\sum_{i=1}^n p_{i,j}}$$

No, the result won't stay positive.



1(e)

Advantages:

- 1) In EM-algorithm the proposed parameters values are always valid for example, probability masses between  $[0,1]$  sums to 1, which is not in the cases of gradient descent.
- 2) In EM-algorithm we don't have to calculate the likelihood to insure it has increased at every step which is not in the case while gradient descent.
- 3) EM method exploits structure of the objective and the variable involved in a manner that they are largely decoupled which allows good convergence rate than gradient descent.