**Department of Computer Science and Technology**

# Machine Learning

Homework 3

## Sahand Sabour

2020280401

# 1 Clustering: Mixture of Multinomials (2 points)

## 1.1 MLE for multinomial (1 point)

The likelihood function for this multinomial distribution is given as

$$P(x|\mu) = \frac{n!}{\prod_i x_i!} \prod_i \mu_i^{x_i}, \quad i = 1, ..., d \tag{1}$$

Taking log from both side of the above equation gives the log-likelihood function

$$\mathcal{L}(\mu) = log(P(x|\mu)) = log(n!) - log(\prod_i x_i!) + log(\prod_i \mu_i^{x_i}) \tag{2}$$

This can be considered a Lagrange problem with the constraint $\sum_i \mu_i = 1$. Hence, the Lagrangian equation can be formulated as

$$\mathcal{L}(\mu) = log(n!) - log(\prod_i x_i!) + log(\prod_i \mu_i^{x_i}) - \lambda(\sum_i \mu_i - 1) \tag{3}$$

where $\lambda$ is Lagrangian multiplier, giving

$$\mathcal{L}(\mu) = log(n!) - \sum_i log(x_i!) + \sum_i x_i log(\mu_i) - \lambda(\sum_i \mu_i - 1) \tag{4}$$

Taking the derivative of the equation with respect to $\mu_i$ and setting it to 0 gives

$$\frac{\partial \mathcal{L}}{\partial \mu_i} = \frac{\sum_i x_i}{\sum_i \mu_i} - \lambda = 0 \tag{5}$$

Hence, we get that

$$\lambda = \frac{\sum_i x_i}{\sum_i \mu_i} = \frac{n}{1} = n \tag{6}$$

Accordingly, we could derive the maximum-likelihood estimator $\mu_i$ as

$$\mu_i = \frac{x_i}{\lambda} = \frac{x_i}{n}, \quad i = 1, ..., d \tag{7}$$

## 1.2 EM for mixture of multinomials (1 point)

Initially, for each document $d$, a latent topic $c_d$ is generated as follows:

$$p(c_d = k) = \pi_k \quad where \quad k = 1, 2, ..., K \tag{8}$$

Accordingly, given a topic $\mu_k$, the bag of words for d is generated:

$$p(d|c_d = k) = \frac{n_d!}{\prod_w T_{dw}!} \prod_w \mu_{w_k}^{T_{dw}} \quad where \quad n_d = \sum_w T_{dw} \tag{9}$$

Combining the above two equations gives

$$p(d) = \sum_{k=1}^{K} p(d|c_d = k)p(c_d = k)$$

$$p(d) = \frac{n_d!}{\prod_w T_{dw}!} \sum_{k=1}^{K} \pi_k \prod_w \mu_{w_k}^{T_{dw}}$$

**E-Step:** estimate the responsibilities.

$$\gamma(c_{dk}) = \frac{\pi_k \mathcal{N}(d_n|\mu_k, \sum_k)}{\sum_j \pi_j \mathcal{N}(d_n|\mu_j, \sum_j)} \tag{10}$$

**M-Step:** re-estimate the parameters. Initially, we have the log likelihood

$$log\, p(\mathcal{D}|\theta) = \sum_{k=1}^{K} log(\sum_j \pi_j \mathcal{N}(d_j \mu_j)) \tag{11}$$

Accordingly, we need to consider this log likelihood as a Lagrangian optimization problem and solve it with respect to $\pi_k$ and $\mu_k$.

$$L = \sum_{k=1}^{K} log(\sum_j \pi_j \mathcal{N}(d_j \mu_j)) - \sum_j \lambda(\mu_j - 1) \tag{12}$$

$$\frac{\delta L}{\delta \pi_k} \tag{13}$$

# 2 PCA (2 points)

## 2.1 Minimum Error Formulation (2 points)

Assuming that we have a set of complete orthonormal basis $\{\mu_i\}$, where $i \in [1, p]$, we have that $\mu_i^T \mu_j = \delta_{ij}$ and each data point can be represented as $x_n = \sum_i a_{ni}\mu_i$. Accordingly, due to orthonormal property, we can get that

$$a_{ni} = x_n^T \mu_i \tag{14}$$

Inserting this in the data point representation gives

$$x_n = \sum_i (x_n^T \mu_i)\mu_i \tag{15}$$

For this approach, the aim is to formulate PCA as minimizing the mean-squared-error of a low-dimensional approximation of the given basis. Hence, we assume a low-dimensional approximation of the point representation as follows

$$\widetilde{x}_n = \sum_i^d z_{ni} + \sum_{i=d+1}^p b_i \mu_i \quad \text{where } b_i \text{ are constants for all points} \tag{16}$$

Therefore, the best approximation is to minimize the following error

$$\min_{U,z,b} J := \frac{1}{N} \sum_{n=1}^N ||x_n - \widetilde{x}_n||^2 \tag{17}$$

Consequently, we have that

$$
\begin{aligned}
J &= \frac{1}{N} \sum_{n=1}^N ||x_n - \widetilde{x}_n||^2 \\
&= \frac{1}{N} \sum_{n=1}^N (x_n - \widetilde{x}_n)^T (x_n - \widetilde{x}_n) \\
&= \frac{1}{N} \sum_{n=1}^N x_n^T x_n - 2x_n^T \widetilde{x}_n + \widetilde{x}_n^T \widetilde{x}_n
\end{aligned}
$$

Inserting equation 10 in the above equation and replacing $\widetilde{x}_n$ gives

$$
\begin{aligned}
J &= \frac{1}{N} \sum_{n=1}^N x_n^T x_n - 2x_n^T \left( \sum_i^d z_{ni}\mu_i + \sum_{i=d+1}^p b_i\mu_i \right) \\
&+ \left( \sum_i^d z_{ni}\mu_i^T + \sum_{i=d+1}^p b_i\mu_i^T \right) \left( \sum_i^d z_{ni}\mu_i + \sum_{i=d+1}^p b_i\mu_i \right)
\end{aligned}
$$

Accordingly, for minimizing this error, we calculate the derivative with respect to $z$ and $b$ and set it to 0.

$$\frac{\delta J}{\delta z_{nj}} = \frac{1}{n}[-2x_n^T\mu_j + \mu_j^T(\sum_i^d z_{ni}\mu_i + \sum_{i=d+1}^p b_i\mu_i) + (\sum_i^d z_{ni}\mu_i^T + \sum_{i=d+1}^p b_i\mu_i^T)\mu_j] = 0$$

$$\frac{\delta J}{\delta z_{nj}} = \frac{1}{n}[-2x_n^T\mu_j + 2\mu_j^T(\sum_i^d z_{ni}\mu_i + \sum_{i=d+1}^p b_i\mu_i)] = 0$$

$$2\mu_j^T(\sum_i^d z_{ni}\mu_i + \sum_{i=d+1}^p b_i\mu_i) = 2x_n^T\mu_j$$

$$\sum_i^d z_{ni}\mu_j^T\mu_i + \sum_{i=d+1}^p b_i\mu_j^T\mu_i = x_n^T\mu_j$$

$$\sum_i^d z_{ni}\delta ij + \sum_{i=d+1}^p b_i\delta ij = z_{ni} + 0 = x_n^T\mu_i$$

Giving $z_{ni} = x_n^T\mu_i$ for $i \in [1, d]$. Similarly, we the derivative with respect to b

$$\frac{\delta J}{\delta b_j} = \frac{1}{n}\sum[-2x_n^T\mu_j + \mu_j^T(\sum_i^d z_{ni}\mu_i + \sum_{i=d+1}^p b_j\mu_i) + (\sum_i^d z_{ni}\mu_i^T + \sum_{i=d+1}^p b_j\mu_i^T)\mu_j] = 0$$

$$\frac{\delta J}{\delta b_j} = \frac{1}{n}\sum[-2x_n^T\mu_j + 2\mu_j^T(\sum_i^d z_{ni}\mu_i + \sum_{i=d+1}^p b_j\mu_i)] = 0$$

$$\sum(\sum_i^d z_{ni}\mu_j^T\mu_i + \sum_{i=d+1}^p b_j\mu_j^T\mu_i) = \sum x_n^T\mu_j$$

$$\sum b_j = Nb_j = \sum x_n^T\mu_j \quad \text{giving} \quad b_j = \sum\frac{1}{n}x_n^T\mu_j = \bar{x}^Tu_j$$

Which in turn gives $b_i = \bar{x}^Tu_i$ for $i \in [d+1, p]$. Accordingly, from equation 9, we can get the displacement lines in the orthogonal subspace as follows

$$x_n - \widetilde{x}_n = \sum_{i=d+1}^p \{(x_n - \bar{x})^T\mu_i\}\mu_i \tag{18}$$

Which produces the following optimization problem for error J

$$\min_{\mu_j} J \quad \text{where} \quad \mu_i^T\mu_i = 1 \tag{19}$$

4

Assuming d=1 (1-dimensional subspace) and p=2 (2-dimensional space), the optimization problem becomes

$$\min_{\mu_2} J = \mu_2^T S \mu_2 \quad \text{where} \quad \mu_2^T \mu_2 = 1 \tag{20}$$

Which gives $S\mu_2 = \lambda_2 \mu_2$, meaning that $\mu_2$ should be chosen as the eigenvector that corresponds to the smaller eigenvalue. Accordingly, the principal subspace is chosen by the eigenvector of the larger eigenvalue.

# 3 Deep Generative Models: Class-conditioned VAE (5 Points)

In the MNIST dataset, there are 10 possible labels for the samples (0-9). Binarizing the labels with the one-hot encoding method, gives a sequence of 10 digits with one 1 and nine 0s. Hence, there could be 10 locations for the 1; the probability of a label l to be one of the 10 labels L would be $p(l = L) = \frac{1}{10} = 0.1$. According to the lecture notes, the variational lower bound for the normal case of VAE was obtained as follows:

$$L(\theta, x) = E_{q(z|x)}[log p(z, x; \theta) - log q(z|x)] = E_{q(z|x)}[log p(x|z; \theta)] - KL(q(z|x)||p(z; \theta))$$

However, it can be noticed that the output of this equation is only dependent on the latent variable z and therefore, does not produce any specific results, which is not practical for our case. Hence, we should modify the lower bound to include the label l of the sample we would like to generate likewise.

$$L(\theta, x, l) = E_{q(z|x,l)}[log p(x, l|z; \theta)] - KL(q(z|x, l)||p(z; \theta))$$

Since $z \sim \mathcal{N}(0, 1)$ for Gaussian, the KL-divergence is as follows:

$$-KL(q(z|x, l)||p(z; \theta)) = \frac{1}{2}(1 + log\sigma^2 - \mu^2 - \sigma^2) \tag{21}$$

Consequently, the expected log-likelihood would be

$$E_{q(z|x,l)}[log p(x, l|z; \theta)] = E_{q(z|x,l)}[-\sum_j \frac{1}{2} log\sigma_j^2 + \frac{(x_{ij} - \mu_{xi})^2}{\sigma^2}] \tag{22}$$

Approximating the above equation with Monte Carlo methods gives

$$E_{q(z|x,l)}[log p(x, l|z; \theta)] \approx \frac{1}{L}\sum_k log p(x, l|z^{(k)}) \quad \text{where} \quad z^{(k)} \sim q(z|x, l) \tag{23}$$

where $z^{(k)}$ is a random variable, which cannot be used for back-propagation. Hence, by utilizing re-parameterization techniques, we have $z^{(k)} = \mu(x, l) + \sigma(x, l).\epsilon^{(k)} = g(x, l, \epsilon^{(k)})$, where g is a deep neural network. The lower bound becomes

$$L(\theta, x, l) = E_{p(\epsilon)}[log\frac{p(g(x, l, \epsilon), x; \theta)}{q(g(x, l, \epsilon)|x; \theta)}] - KL(q(z|x, l)||p(z; \theta))$$

$$L(\theta, x, l) = \frac{1}{L}\sum_{k} logp(x, l|z^{(k)}) + \frac{1}{2}\sum_{i=1}^{j}[1 + log\sigma^2 - \mu^2 - \sigma^2]$$