# ML Homework #3

Yoke Kai Wen, 2020280598

December 27, 2020

# 1 Clustering: Mixture of Multinomials

### 1.1 MLE for multinomial

$$P(\boldsymbol{x}|\boldsymbol{\mu}) = \frac{n!}{\prod_{i} x_{i}!} \prod_{i} \mu_{i}^{x_{i}} = n! \prod_{i} \frac{\mu_{i}^{x_{i}}}{x_{i}!}$$
(1)

Log-likelihood:

$$\begin{split} l(\pmb{\mu}) &= log(P(\pmb{x}|\pmb{\mu})) \\ &= log(n! \prod_{i} \frac{\mu_{i}^{x_{i}}}{x_{i}!}) \\ &= log(n!) + log(\prod_{i} \frac{\mu_{i}^{x_{i}}}{x_{i}!}) \\ &= log(n!) + \sum_{i=1}^{d} log(\frac{\mu_{i}^{x_{i}}}{x_{i}!}) \\ &= log(n!) + \sum_{i=1}^{d} x_{i} log(\mu_{i}) - \sum_{i=1}^{d} x_{i}! \end{split} \tag{2}$$

Form Lagrangian with constraint  $\sum_{i} \mu_{i} = 1$ :

$$L(\boldsymbol{\mu}, \lambda) = l(\boldsymbol{\mu}) + (1 - \sum_{i=1}^{d} \mu_i)$$
(3)

Find  $argmax_{\mu}L(\mu, \lambda)$ 

$$\frac{\partial}{\partial \mu_i} L(\boldsymbol{\mu}, \lambda) = \frac{\partial}{\partial \mu_i} l(\boldsymbol{\mu}) + \frac{\partial}{\partial \mu_i} (1 - \sum_{i=1}^d \mu_i) = 0$$

$$\frac{x_i}{\mu_i} + \lambda(-1) = 0$$

$$\mu_i = \frac{x_i}{\lambda}$$

$$\sum_{i=1}^d \mu_i = \sum_{i=1}^d \frac{x_i}{\lambda}$$
(4)

Result for  $\mu_i$  can be further simplified:

$$LHS = 1, \quad RHS = \frac{1}{\lambda} \sum_{i=1}^{d} x_i = \frac{n}{\lambda} \quad (given)$$

$$1 = \frac{n}{\lambda} \longrightarrow \lambda = n$$

$$\mu_i = \frac{x_i}{n}$$
(5)

#### 1.2 EM for mixture of multinomials

Let  $x_d$  represent 1xW feature vector of document d. Rewriting probability distribution:

$$P(\boldsymbol{x}_{d}|c_{d} = k) = P(\boldsymbol{x}_{d}|\boldsymbol{\mu}_{k}) = n_{d}! \prod_{\omega=1}^{W} \frac{\mu_{\omega k}^{T_{d\omega}}}{T_{d\omega}!}$$

$$\sum_{\omega=1}^{W} T_{d\omega} = n_{d}$$

$$\sum_{\omega=1}^{W} \mu_{\omega k} = 1, \quad \forall k$$

$$\sum_{k=1}^{K} \pi_{k} = 1$$

$$(6)$$

Expectation step is similar to GMM, we also introduce a q distribution for latent variable:

$$q(c_{d} = k) = P(c_{d} = k | \boldsymbol{x}_{d}, \boldsymbol{\pi}, \boldsymbol{\mu}) = \frac{\pi_{k} P(\boldsymbol{x}_{d} | \boldsymbol{\mu}_{k})}{\sum_{j=1}^{k} \pi_{j} P(\boldsymbol{x}_{d} | \boldsymbol{\mu}_{j})} = \frac{\pi_{k} n_{d}! \prod_{\omega=1}^{W} \frac{\mu_{\omega k}^{T_{d\omega}}}{T_{d\omega}!}}{\sum_{j=1}^{K} \pi_{j} n_{d}! \prod_{\omega=1}^{W} \frac{\mu_{\omega j}^{T_{d\omega}}}{T_{d\omega}!}}$$

$$= \frac{\pi_{k} \prod_{\omega=1}^{W} \mu_{\omega k}^{T_{d\omega}}}{\sum_{i=1}^{K} \pi_{j} \prod_{\omega=1}^{W} \mu_{\omega j}^{T_{d\omega}}}$$

$$= \frac{\pi_{k} \prod_{\omega=1}^{W} \mu_{\omega k}^{T_{d\omega}}}{\sum_{i=1}^{K} \pi_{j} \prod_{\omega=1}^{W} \mu_{\omega j}^{T_{d\omega}}}$$
(7)

EM algorithm defines lower bound for log-likelihood by applying Jensen's inequality.

$$log(P(\mathbf{x}_{d}|\mathbf{\pi}.\boldsymbol{\mu}) = log(\sum_{k=1}^{K} q(c_{d} = k) \frac{P(\mathbf{x}_{d}, c_{d} = k|\pi_{k}, \boldsymbol{\mu}_{k})}{q(c_{d} = k)}) \ge \sum_{k=1}^{K} q(c_{d} = k) log(\frac{P(\mathbf{x}_{d}, c_{d} = k|\pi_{k}, \boldsymbol{\mu}_{k})}{q(c_{d} = k)})$$

$$= L_{i}(\boldsymbol{\pi}, \boldsymbol{\mu}.\boldsymbol{q})$$
(8)

Maximisation step with respect to  $\mu, \pi$ :

First, include all documents in the lower bound for the log-likelihood:

$$L(\boldsymbol{\pi}, \boldsymbol{\mu}.\boldsymbol{q}) = \sum_{d=1}^{D} L_{i}(\boldsymbol{\pi}, \boldsymbol{\mu}.\boldsymbol{q})$$

$$= \sum_{d=q}^{D} \sum_{k=1}^{K} q(c_{d} = k) log(\frac{P(\boldsymbol{x}_{d}, c_{d} = k | \pi_{k}, \boldsymbol{\mu}_{k})}{q(c_{d} = k)})$$

$$= \sum_{d=1}^{D} \sum_{k=1}^{K} q(c_{d} = k) log(P(\boldsymbol{x}_{d}, c_{d} = k) P(\boldsymbol{x}_{d} | c_{d} = k, \pi_{k}, \boldsymbol{\mu}_{k})) - \sum_{d=1}^{D} \sum_{k=1}^{K} q(c_{d} = k) log(q(c_{d} = k))$$
(9)

Last term in  $L(\boldsymbol{\pi}, \boldsymbol{\mu}.\boldsymbol{q})$  can be dropped during maximisation because it does not contain  $\boldsymbol{\pi}, \boldsymbol{\mu}$  terms, we can call it *const*. Then formulate Lagrangian optimisation problem from  $L(\boldsymbol{\pi}, \boldsymbol{\mu}.\boldsymbol{q})$  with constraints  $\sum_{\omega=1}^{W} \mu_{\omega k} = 1, \forall k$  and  $\sum_{k=1}^{K} \pi_k = 1$ .

$$L' = \sum_{d=1}^{D} \sum_{k=1}^{K} k = 1^{K} q(c_d = k) log(P(c_d = k) | \pi_k, \boldsymbol{\mu}_k) - \lambda_1([\sum_{\omega=1}^{W} \mu_{\omega k}] - 1) - \lambda_2([\sum_{k=1}^{K} \pi_k] - 1) + const$$
(10)

Then we differentiate and set differential to zero, first with respect to  $\mu_{\omega k}$ :  $\frac{\partial L'}{\partial \mu_{\omega k}} = 0$ .

$$\frac{\partial}{\partial \mu_{\omega k}} \left[ \left( \sum_{d=1}^{D} q(c_d = k) log(\pi_k n_d! \prod_{\omega=1}^{W} \frac{\mu_{\omega k}^{T_{d\omega}}}{T_{d\omega}!}) \right) - \lambda_1 \mu_{\omega k} \right] = 0$$

$$\frac{\partial}{\partial \mu_{\omega k}} \left[ \left( \sum_{d=1}^{D} q(c_d = k) log(\pi_k n_d!) \right) + \left( \sum_{d=1}^{D} q(c_d = k) \sum_{\omega=1}^{W} log(\frac{\mu_{\omega k}^{T_{d\omega}}}{T_{d\omega}!}) \right) - \lambda_1 \mu_{\omega k} \right] = 0$$

$$\frac{\partial}{\partial \mu_{\omega k}} \left[ \left( \sum_{d=1}^{D} q(c_d = k) \sum_{\omega=1}^{W} \left[ T_{d\omega} log(\mu_{\omega k}) - log(T_{d\omega}!) \right] \right) - \lambda_1 \mu \right] = 0$$

$$\frac{\partial}{\partial \mu_{\omega k}} \left[ \left( \sum_{d=1}^{D} q(c_d = k) T_{d\omega} log(\mu_{\omega k}) \right) - \lambda_1 \mu_{\omega k} \right] = 0$$

$$\frac{\sum_{d=1}^{D} q(c_d = k) T_{d\omega}}{\mu_{\omega k}} - \lambda_1 = 0$$

$$\mu_{\omega k} = \frac{\sum_{d=1}^{D} q(c_d = k) T_{d\omega}}{\lambda_1}$$
(11)

Simplify expression for  $\mu_{\omega k}$ :

$$\sum_{\omega=1}^{W} \mu_{\omega k} = \frac{\sum_{\omega=1}^{W} \sum_{d=1}^{D} q(c_d = k) T_{d\omega}}{\lambda_1}$$

$$LHS = 1, \quad \sum_{\omega=1}^{W} T_{d\omega} = n_d \quad (given)$$

$$\lambda_1 = \sum_{d=1}^{D} q(c_d = k) n_d$$

$$\mu_{\omega k} = \frac{\sum_{d=1}^{D} q(c_d = k) T_{d\omega}}{\sum_{d=1}^{D} q(c_d = k) n_d}$$
(12)

Now, we differentiate with respect to  $\pi_k$ , then  $\frac{\partial L'}{\partial \pi_k} = 0$ .

$$\frac{\partial}{\partial \pi_k} \left[ \left( \sum_{d=1}^D q(c_d = k) log(\pi_k n_d! \prod_{\omega=1}^W \frac{\mu_{\omega k}^{T_{d\omega}}}{T_{d\omega}!} \right) \right) - \lambda_2 \pi_k \right] = 0$$

$$\frac{\partial}{\partial \pi_k} \left[ \left( \sum_{d=1}^D q(c_d = k) log(\pi_k) \right) + \left( \sum_{d=1}^D q(c_d = k) \sum_{\omega=1}^W log(n_d! \frac{\mu_{\omega k}^{T_{d\omega}}}{T_{d\omega}!} \right) \right) - \lambda_2 \pi_k \right] = 0$$

$$\frac{\partial}{\partial \pi_k} \left[ \sum_{d=1}^D q(c_d = k) log(\pi_k) - \lambda_2 \pi_k \right] = 0$$

$$\frac{\sum_{d=1}^D q(c_d = k)}{\pi_k} - \lambda_2 = 0$$

$$\pi_k = \frac{\sum_{d=1}^D q(c_d = k)}{\lambda_2}$$
(13)

Simplify expression for  $\pi_k$ :

$$\sum_{k=1}^{K} \pi_k = \frac{\sum_{k=1}^{K} \sum_{d=1}^{D} q(c_d = k)}{\lambda_2}$$

$$LHS = 1, \quad \sum_{k=1}^{K} q(c_d = k) = 1 \quad (given)$$

$$\lambda_2 = \sum_{d=1}^{D} 1 = D$$

$$\pi_k = \frac{\sum_{d=1}^{D} q(c_d = k)}{D}$$
(14)

Therefore, the update rules are:

For the Expectation step:

$$q(c_d = k) = \frac{\pi_k \prod_{\omega=1}^{W} \mu_{\omega k}^{T_{d\omega}}}{\sum_{j=1}^{K} \pi_j \prod_{\omega=1}^{W} \mu_{\omega j}^{T_{d\omega}}}$$

For the Maximisation step:

$$\mu_{\omega k} = \frac{\sum_{d=1}^{D} q(c_d = k) T_{d\omega}}{\sum_{d=1}^{D} q(c_d = k) n_d}$$
$$\pi_k = \frac{\sum_{d=1}^{D} q(c_d = k)}{D}$$

## 2 PCA

### 2.1 Minimum Error Formulation

Task: Given that actual data point

$$x_n = \sum_{i=1}^{P} (x_n^T \mu_i) \mu_i$$

and low dimensional approximation

$$\widetilde{x_n} = \sum_{i=1}^d z_{ni} \mu_i + \sum_{i=d+1}^P b_i \mu_i$$

and error

$$J = \frac{1}{N} \sum_{n=1}^{N} ||x_n - \widetilde{x_n}||^2$$

prove that the best approximation to minimise the error J is:  $z_{ni} = x_n^T \mu_i$  for  $i=1..d, b_i = \overline{x_n}^T \mu_i$  for i=(d+1)..p

First, reformulate J to simplify the vector norm so that differentiation later will be easier.

$$J = \frac{1}{N} \sum_{n=1}^{N} (x_n - \widetilde{x_n})^T (x_n - \widetilde{x_n})$$

$$= \frac{1}{N} \sum_{n=1}^{N} (x_n^T x_n - 2x_n^T \widetilde{x_n} + \widetilde{x_n}^T \widetilde{x_n})$$
(15)

Simplify expressions for  $x_n^T x_n$ ,  $x_n^T \widetilde{x_n}$  and  $\widetilde{x_n}^T \widetilde{x_n}$  by using orthonormal property of  $\mu_i$ , i.e.

$$\mu_i^T \mu_j = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases}$$
 (16)

$$x_n^T x_n = (\sum_{i=1}^P (x_n^T \mu_i) \mu_i)^T (\sum_{i=1}^P (x_n^T \mu_i) \mu_i)$$

$$= \sum_{i=1}^P (x_n^T \mu_i)^2$$
(17)

$$x_n^T \widetilde{x_n} = (\sum_{i=1}^P (x_n^T \mu_i) \mu_i)^T (\sum_{i=1}^d z_{ni} \mu_i + \sum_{i=d+1}^P b_i \mu_i)$$

$$= \sum_{i=1}^d (x_n^T \mu_i) z_{ni} + \sum_{i=d+1}^P (x_n^T \mu_i) b_i$$
(18)

$$\widetilde{x_n}^T \widetilde{x_n} = (\sum_{i=1}^d z_{ni} \mu_i + \sum_{i=d+1}^P b_i \mu_i)^T (\sum_{i=1}^d z_{ni} \mu_i + \sum_{i=d+1}^P b_i \mu_i)$$

$$= \sum_{i=1}^d z_{ni}^2 + \sum_{i=d+1}^P b_i^2$$
(19)

Then, J can be rewritten as:

$$J = \frac{1}{N} \sum_{n=1}^{N} \left( \sum_{i=1}^{P} (x_n^T \mu_i)^2 - 2 \sum_{i=1}^{d} (x_n^T \mu_i) z_{ni} - 2 \sum_{i=d+1}^{P} (x_n^T \mu_i) b_i + \sum_{i=1}^{d} z_{ni}^2 + \sum_{i=d+1}^{P} b_i^2 \right)$$
(20)

Then, differentiate J with respect to  $z_{ni}$  and  $b_i$  to get the best approximation to minimise J. Differentiating with respect to  $z_{ni}$ ,  $\frac{\partial J}{\partial z_{ni}} = 0$ :

$$\frac{1}{N}(-2x_n^T \mu_i + 2z_{ni}) = 0$$

$$z_{ni} = x_n^T \mu_i \quad (proven)$$
(21)

Differentiating with respect to  $b_i$ ,  $\frac{\partial J}{\partial b_i} = 0$ :

$$\frac{1}{N}(-2x_n^T\mu_i + 2b_i) = 0$$

$$\sum_{n=1}^N b_i = \sum_{n=1}^N x_n^T\mu_i$$

$$Nb_i = \sum_{n=1}^N x_n^T\mu_i$$

$$b_i = \frac{1}{N} \sum_{n=1}^N x_n^T\mu_i = \overline{x_n}^T\mu_i \quad (proven)$$
(22)

# 3 Deep Generative Models: Class-conditioned VAE

### 3.1 Derive algorithm for CVAE

Algorithm for CVAE is same as VAE, except with class condition y. MLE learning:

$$max_{\theta}log(P_{\theta}(x|y))$$

But  $P_{\theta}(x|y)$  is conditional marginal likelihood of x, and is hard to compute, since we only have the joint likelihood  $P_{\theta}(x,z|y) = P_{\theta}(x|z,y)P(z|y)$ , but computing  $P_{\theta}(x|y) = \int P_{\theta}(x,z|y)dz$  is intractable.

So instead of maximising  $log(P_{\theta}(x|y))$  directly, we choose to optimise the variational lower bound of it, and introduce the class conditioned variational posterior  $q_{\phi}(z|x,y)$ .

$$log(P_{\theta}(x|y)) \ge log(P_{\theta}(x|y)) - KL(q_{\phi}(z|x,y)||P_{\theta}(z|x,y))$$

$$= E_{q_{\phi}(z|x,y)}[log(P_{\theta}(x,z|y)) - log(q_{\phi}(z|x,y))]$$

$$= L(\theta,\phi)$$
(23)

#### 3.2 Implement CVAE with ZhuSuan on MNIST dataset

I modified the vae.py example code in Zhusuan to include the class condition. The CVAE python script cvae\_mnist2.py is in the attached zip file. The main modifications to the example code is the concatenation of the class variable y with the latent variable z in the generator (decoder) and concatenation of the class variable y with the observed variable x in the encoder as seen in the code fragments below.

```
1 @zs.meta_bayesian_net(scope="gen", reuse_variables=True)
def build_gen(y, x_dim, z_dim, y_dim, n):
      bn = zs.BayesianNet()
      # z ~ N(0,I)
      z_mean = tf.zeros([n, z_dim])
      z = bn.normal("z", z_mean, std=1., group_ndims=1)
6
      \# Concatenate z and y
8
      z = tf.concat(axis=1, values=[z,y])
9
10
      \# x_{logits} = f_{NN}(z)
      h = tf.layers.dense(z, 500, activation=tf.nn.relu)
      h = tf.layers.dense(h, 500, activation=tf.nn.relu)
      x_logits = tf.layers.dense(h, x_dim)
14
      x_mean = bn.deterministic("x_mean", tf.sigmoid(x_logits))
16
17
      # x ~ Bernoulli(x_logits)
      bn.bernoulli("x", x_logits, group_ndims=1, dtype=tf.float32)
19
      return bn
20
21
22
23 @zs.reuse_variables(scope="q_net")
24 def build_q_net(x, y, z_dim, y_dim):
      bn = zs.BayesianNet()
      # concatenate x and y
26
      x = tf.concat(axis=1, values=[x,y])
27
      h = tf.layers.dense(x, 500, activation=tf.nn.relu)
      h = tf.layers.dense(h, 500, activation=tf.nn.relu)
29
30
      z_mean = tf.layers.dense(h, z_dim)
      z_logstd = tf.layers.dense(h, z_dim)
      bn.normal("z", z_mean, logstd=z_logstd, group_ndims=1)
    return bn
```

#### 3.3 Visualise CVAE generations

The CVAE model trained on MNIST produced the following visualisations at the first epoch (Figure 1) and at the 30th epoch (Figure 2).



Figure 1: CVAE generations in first epoch



Figure 2: CVAE generations in 30th epoch