

## Module 4F10: STATISTICAL PATTERN RECOGNITION

## Solutions to Examples Paper 1

1. Average risk in choosing class  $\omega_i$  is

$$\begin{aligned} R(\omega_i|\mathbf{x}) &= \sum_{j=1}^c \lambda(\omega_i|\omega_j)P(\omega_j|\mathbf{x}) \\ &= 0.P(\omega_i|\mathbf{x}) + \sum_{j=1, j \neq i}^c \lambda_s P(\omega_j|\mathbf{x}) \end{aligned}$$

where  $\lambda(\omega_i|\omega_j)$  is used to mean the cost of choosing class  $\omega_i$  where the true class is  $\omega_j$ .

Hence

$$R(\omega_i|\mathbf{x}) = \lambda_s (1 - P(\omega_i|\mathbf{x}))$$

Associate  $\mathbf{x}$  with class  $\omega_i$  if highest posterior class probability and the average risk is less than the cost of rejection

$$\begin{aligned} \lambda_s (1 - P(\omega_i|\mathbf{x})) &\leq \lambda_r \\ P(\omega_i|\mathbf{x}) &\geq 1 - \lambda_r/\lambda_s \end{aligned}$$

If the ratio  $\lambda_r/\lambda_s$  is close to 1 then the reject region will tend to zero. If  $\lambda_r/\lambda_s$  is close to zero then nearly all examples will be rejected.

2. The computational cost for the 3 systems are

- Diagonal covariance the cost is  $2d$  multiply accumulates.
- Full covariance the cost is  $d^2 + d$  multiply accumulates.
- $M$  component diagonal system the cost is  $2Md$ .

In all cases the inverse covariance matrix is stored and the portion of the Gaussian PDF not dependent on the observations is stored as a constant.

As a practical matter for the Gaussian mixture case,  $\log(P(\omega_m))$  is added to the constant in advance and the computation can be arranged to use in the sum (in the log domain)

$$\log(\exp(l_i(\mathbf{x})) + \exp(l_j(\mathbf{x}))) = l_i(\mathbf{x}) + \log(1 + \exp(l_j(\mathbf{x}) - l_i(\mathbf{x})))$$

assuming that  $l_i(\mathbf{x}) \geq l_j(\mathbf{x})$  and

$$l_i(\mathbf{x}) = \log(P(\omega_i)) + \log(\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i))$$

This often saves one exponential calculation, the exponential is only calculated when the value of the difference in logs is above a threshold, and handles any dynamic range issues.

Note an approximation that is sometimes used for a GMM is that the overall likelihood approximated by using only the largest of the component log-likelihoods i.e.

$$\log(p(\mathbf{x})) \approx \max_m (\log(P(\omega_m)) + \log(\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)))$$

This saves adding the logs above but is not useful for training where more precise calculation is required.

3. (a) In this case, the variance is equal to 1 for each of the single dimensional Gaussians. The log likelihood of the model for single dimensional Gaussians and two mixture components is

$$l(\boldsymbol{\theta}) = \sum_{k=1}^n \ln \left[ \sum_{m=1}^2 c_m \frac{1}{(2\pi)^{1/2}} \exp \left\{ \frac{-(x_k - \mu_m)^2}{2} \right\} \right]$$

Here  $n = 9$  and the question asks for the likelihood with  $c_1 = c_2 = 0.5$  which can be calculated as

$$\prod_{k=1}^9 \sum_{m=1}^2 \frac{1}{q} \exp \left\{ \frac{-(x_k - \mu_m)^2}{2} \right\}$$

where  $1/q = 0.199$ . Substituting in the data values and the mean values yields the total likelihood of the data as  $2.262 \times 10^{-7}$ .

(b) To compute the re-estimated means and component priors / mixture weights, find the posterior probabilities of each mixture component for each data sample and accumulate numerator and denominator statistics. This is most easily done by writing a small program/script. The posterior probabilities of each mixture component is given in the table below.

$x$	$P(\text{comp1} x)$	$P(\text{comp2} x)$
-1.5	0.9933	0.0067
-0.5	0.9526	0.0474
0.1	0.8581	0.1419
0.3	0.8022	0.1978
0.9	0.5498	0.4502
1.3	0.3543	0.6457
1.9	0.1419	0.8581
2.3	0.0691	0.9309
3.0	0.0180	0.9820

Then for each mixture component accumulate for mean re-estimation the numerator  $\sum_{k=1}^9 x_k P(m|x_k)$  and for the denominator  $\sum_{k=1}^9 P(m|x_k)$ . The same statistics used for the denominator are also needed for component prior re-estimation.

Computing these values leads to  $\hat{\mu}_1 = -0.0426$  ;  $\hat{\mu}_2 = 1.878$  ;  $\hat{c}_1 = 0.5266$ ;  $\hat{c}_2 = 0.4734$ . Use of these values to re-compute the likelihood yields an increase over the initial values.

4. From lecture notes we can write the auxiliary function as

$$\mathcal{Q}(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}) = \sum_{m=1}^M \sum_{i=1}^n P(\omega_m|\mathbf{x}_i) \sum_{k=1}^d [x_{ik} \log(\lambda_{mk}) + (1 - x_{ik}) \log(1 - \lambda_{mk})]$$

Differentiate this with respect to  $\lambda_{qr}$  give

$$\frac{\partial \mathcal{Q}(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}})}{\partial \lambda_{qr}} = \sum_{i=1}^n P(\omega_q|\mathbf{x}_i) \left[ \frac{x_{ir}}{\lambda_{qr}} - \frac{(1 - x_{ir})}{(1 - \lambda_{qr})} \right]$$

Equating to zero gives

$$(1 - \lambda_{qr}) \sum_{i=1}^n P(\omega_q|\mathbf{x}_i) x_{ir} = \lambda_{qr} \sum_{i=1}^n P(\omega_q|\mathbf{x}_i) (1 - x_{ir})$$

Rearranging yields the answer.

5. We can write

$$x_i = t_i + z$$

where  $z$  is Gaussian distributed, mean 0 and variance 1.

(a) Since the two are independent and both Gaussian distributed we know that

$$p(x_i|\theta) = \mathcal{N}(x_i; \mu, \sigma^2 + 1)$$

We can write the log-likelihood of the training data as

$$\begin{aligned} l(\boldsymbol{\theta}) &= \sum_{i=1}^n \log(\mathcal{N}(x_i, \mu, 1 + \sigma^2)) \\ &= \sum_{i=1}^n -\frac{1}{2} \left( \log(2\pi(1 + \sigma^2)) + \frac{(x_i - \mu)^2}{1 + \sigma^2} \right) \end{aligned}$$

Differentiating

$$\frac{\partial l(\boldsymbol{\theta})}{\partial \mu} = \sum_{i=1}^n \frac{x_i - \mu}{1 + \sigma^2}$$

Equating to zero gives

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

For the variance

$$\frac{\partial l(\boldsymbol{\theta})}{\partial \sigma^2} = \sum_{i=1}^n -\frac{1}{2} \left( \frac{1}{(1 + \sigma^2)} - \frac{(x_i - \mu)^2}{(1 + \sigma^2)^2} \right)$$

Equating to zero and using the ML estimate for  $\mu$

$$\sigma^2 = \left( \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \right) - 1$$

(b) This is the same set up as described in lecture. Let  $z_i$  be the noise associated with observation  $i$ . So

$$p(x_i|z_i, \theta) = \mathcal{N}(x_i; \mu + z_i, \sigma^2)$$

We first need to compute the posterior  $p(z_i|x_i, \theta)$

$$\begin{aligned} p(z_i|x_i, \theta) &= \frac{p(x_i|z_i, \theta)p(z_i)}{p(x_i|\theta)} \\ &= \mathcal{N}\left(z_i; \frac{(x_i - \mu)}{(1 + \sigma^2)}, \frac{\sigma^2}{(1 + \sigma^2)}\right) \end{aligned}$$

So writing down the auxiliary function

$$\begin{aligned} \mathcal{Q}(\theta, \hat{\theta}) &= \sum_{i=1}^n \int (p(z_i|x_i, \theta) \log(p(x_i, z_i|\hat{\theta}))) dz_i \\ &= \sum_{i=1}^n \int (p(z_i|x_i, \theta) \log(p(x_i|z_i, \hat{\theta}))) dz_i \\ &\quad + \sum_{i=1}^n \int (p(z_i|x_i, \theta) \log(p(z_i))) dz_i \end{aligned}$$

The second term is not dependent on the new model parameters, the distribution of  $z_i$  is known. This leaves the first term. From the previous definitions

$$\begin{aligned} \tilde{\mathcal{Q}}(\theta, \hat{\theta}) &= \sum_{i=1}^n \int p(z_i|x_i, \theta) \log(p(x_i|z_i, \hat{\theta})) dz_i \\ &= \sum_{i=1}^n \int p(z_i|x_i, \theta) \left[ \log\left(\frac{1}{\sqrt{2\pi\hat{\sigma}^2}}\right) - \frac{(x_i - z_i - \hat{\mu})^2}{2\hat{\sigma}^2} \right] dz_i \\ &= \sum_{i=1}^n \left[ \log\left(\frac{1}{\sqrt{2\pi\hat{\sigma}^2}}\right) - \frac{(x_i - \hat{\mu})^2 - 2(x_i - \hat{\mu})\mathcal{E}\{z_i|\theta, x_i\} + \mathcal{E}\{z_i^2|\theta, x_i\}}{2\hat{\sigma}^2} \right] \end{aligned}$$

We know that

$$\begin{aligned}\mathcal{E}\{z_i|\theta, x_i\} &= \frac{(x_i - \mu)}{(1 + \sigma^2)} \\ \mathcal{E}\{z_i^2|\theta, x_i\} &= \frac{\sigma^2}{(1 + \sigma^2)} + \left(\frac{(x_i - \mu)}{(1 + \sigma^2)}\right)^2\end{aligned}$$

Differentiating with respect to  $\hat{\mu}$  gives

$$\frac{\partial \tilde{Q}(\theta, \hat{\theta})}{\partial \hat{\mu}} = \sum_{i=1}^n \frac{1}{\hat{\sigma}^2} (x_i - \hat{\mu} - \mathcal{E}\{z_i|\theta, x_i\})$$

so

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \left( x_i - \frac{(x_i - \mu)}{(1 + \sigma^2)} \right) = \frac{1}{n} \sum_{i=1}^n \frac{(\sigma^2 x_i + \mu)}{(1 + \sigma^2)}$$

Differentiating with respect to  $\sigma^2$

$$\frac{\partial \tilde{Q}(\theta, \hat{\theta})}{\partial \sigma^2} = \frac{1}{2} \sum_{i=1}^n \left[ \frac{-1}{\hat{\sigma}^2} + \frac{(x_i - \hat{\mu})^2 - 2(x_i - \hat{\mu})\mathcal{E}\{z_i|\theta, x_i\} + \mathcal{E}\{z_i^2|\theta, x_i\}}{(\hat{\sigma}^2)^2} \right]$$

Equating to zero gives

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \left[ (x_i - \hat{\mu})^2 - 2(x_i - \hat{\mu})\mathcal{E}\{z_i|\theta, x_i\} + \mathcal{E}\{z_i^2|\theta, x_i\} \right]$$

This problem is simple to solve using standard optimisation techniques. For EM the correct answer is eventually obtained, but only after many iterations. For situations where direct optimisation is not simple, EM is useful, for example mixture models. Gradient descent could be used in these situations, but do not have the guaranteed stability of EM.

6. It is possible to express  $\mathbf{A}$  as

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -1 & 1 \end{bmatrix}$$

These expressions all match to experts. Thus  $\text{diag}()$  yields a matrix with the vector as the leading diagonal)

$$\boldsymbol{\mu} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 2 \\ -1 \end{bmatrix}; \quad \boldsymbol{\Sigma} = \text{diag} \left( \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \right)$$

5

Comparing the expansion to a standard Gaussian yields

$$\begin{aligned}\boldsymbol{\Sigma}^{-1} &= \mathbf{A}'\mathbf{A} \\ \bar{\boldsymbol{\mu}} &= \boldsymbol{\Sigma}\mathbf{A}'\boldsymbol{\mu}\end{aligned}$$

The mean can be used as the trajectory for speech synthesis. In numbers

$$\hat{\mathbf{x}} = \bar{\boldsymbol{\mu}} = \begin{bmatrix} 0.9231 \\ 1.7692 \\ 1.3846 \end{bmatrix}$$

When the two experts describing the “gradient” are removed  $\hat{\mathbf{x}}$  simply becomes the means

$$\hat{\mathbf{x}} = \bar{\boldsymbol{\mu}} = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}$$

7. From the form of the RBM it is possible to write

$$P(\mathbf{h}|\mathbf{x}, \boldsymbol{\theta}) = \prod_{j=1}^J P(h_j|\mathbf{x}, \boldsymbol{\theta})$$

The conditional distribution for dimensions  $j$  can be written as (using Bayes’ rule)

$$P(h_j = 1|\mathbf{x}, \boldsymbol{\theta}) = \frac{P(h_j = 1, \mathbf{x}|\boldsymbol{\theta})}{p(\mathbf{x}|\boldsymbol{\theta})} \propto P(h_j = 1, \mathbf{x}|\boldsymbol{\theta})$$

Substituting  $h_j = 1$  into the energy function yields:

$$P(h_j = 1, \mathbf{x}|\boldsymbol{\theta}) \propto \exp \left( - \sum_{i=1}^d \frac{(x_i - a_i)^2}{2\sigma_i^2} + b_j + \sum_i \frac{x_i}{\sigma_i} w_{ij} \right)$$

For each dimension  $h_j$  is either 1 or zero. It is also possible to write

$$P(h_j = 0, \mathbf{x}|\boldsymbol{\theta}) \propto \exp \left( - \sum_{i=1}^d \frac{(x_i - a_i)^2}{2\sigma_i^2} \right)$$

Combining these conditions together and cancelling terms yields

$$\begin{aligned}P(h_j = 1|\mathbf{x}, \boldsymbol{\theta}) &= \frac{P(h_j = 1, \mathbf{x}|\boldsymbol{\theta})}{P(h_j = 0, \mathbf{x}|\boldsymbol{\theta}) + P(h_j = 1, \mathbf{x}|\boldsymbol{\theta})} = \frac{\exp(b_j + \sum_i \frac{x_i}{\sigma_i} w_{ij})}{1 + \exp(b_j + \sum_i \frac{x_i}{\sigma_i} w_{ij})} \\ &= \frac{1}{1 + \exp(-b_j - \sum_i \frac{x_i}{\sigma_i} w_{ij})}\end{aligned}$$

6

The distribution for  $p(\mathbf{x}|\mathbf{h}, \boldsymbol{\theta})$  can be found in a similar way. Here we could compute  $P(\mathbf{h}|\boldsymbol{\theta})$ .

$$P(\mathbf{h}|\boldsymbol{\theta}) = \int \exp \left( - \left( \sum_{i=1}^d \frac{(x_i - a_i)^2}{2\sigma_i^2} - \sum_{j=1}^J b_j h_j - \sum_{i,j} \frac{x_i}{\sigma_i} h_j w_{ij} \right) \right) d\mathbf{x}$$

but again this is only necessary as a normalisation term. Thus considering only the numerator term

$$p(\mathbf{x}|\mathbf{h}, \boldsymbol{\theta}) \propto p(\mathbf{x}, \mathbf{h}|\boldsymbol{\theta}) = \exp \left( - \sum_{i=1}^d \frac{(x_i - a_i)^2}{2\sigma_i^2} + \sum_{i,j} \frac{x_i}{\sigma_i} h_j w_{ij} \right)$$

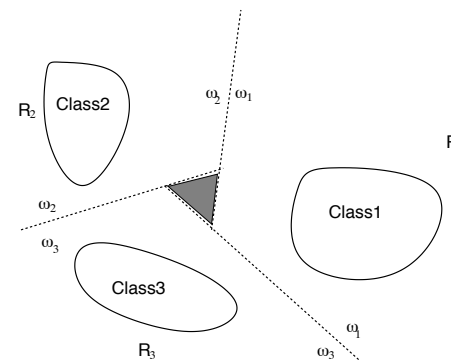
Extracting out the terms for the first and second moments for a dimension  $x_i$

$$-\frac{x_i^2}{2\sigma_i^2} + \frac{x_i}{\sigma_i^2} \left( a_i + \sigma_i \sum_j h_j w_{ij} \right)$$

Comparing to the standard Gaussian this directly yields

$$p(x_i|\mathbf{h}, \boldsymbol{\theta}) = \mathcal{N}(x_i; a_i + \sigma_i \sum_j h_j w_{ij}, \sigma_i^2)$$

8. For the one v one classifiers: it is necessary to train and decode with  $K(K-1)/2$  classifiers. The no decision region is shown below



The region marked in gray is “no decision”. In this case it appears to be both class 1 and class 2

For the one v rest classifiers: it is necessary to train and decode with  $(K-1)$  classifiers. The no decision region is shown below

