

## Practice Set 3 (Problem 1)

$$X = \{x_1, x_2, \dots, x_N\}$$

$$Y = \{y_1, y_2, \dots, y_N\}$$

The likelihood (both  $X$  and  $Y$  are modeled here)

$$P(X, Y | \Theta) = \prod_{n=1}^N P(x_n, y_n | \Theta)$$

Indicator  
 $\uparrow$   
 $\mathbb{I}[y_n=k]$

$$\Theta = \{\mu_k, \Sigma_k, \pi_k\}_{k=1}^K = \prod_{n=1}^N \prod_{k=1}^K \left[ P(x_n, y_n=k | \Theta) \right]$$

$\underbrace{\quad \quad \quad}_{\text{Note that only one term in the product will be selected based on the value of } y_n}$

Thus log-likelihood will be

$$\begin{aligned} \log P(X, Y | \Theta) &= \sum_{n=1}^N \sum_{k=1}^K \mathbb{I}[y_n=k] \left[ \log P(x_n, y_n=k | \Theta) \right] \\ &= \sum_{n=1}^N \sum_{k=1}^K \mathbb{I}[y_n=k] \left[ \log P(x_n | y_n=k, \Theta) P(y_n=k | \Theta) \right] \end{aligned}$$

$\downarrow N(x_n | \mu_k, \Sigma_k)$        $\downarrow \pi_k$

$$L(\Theta) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{I}[y_n=k] \left[ \log N(x_n | \mu_k, \Sigma_k) + \log \pi_k \right]$$

To optimize w.r.t.  $\pi_k$ , we take partial derivatives w.r.t.  $\pi_k$  BUT need to use the constraint  $\sum_k \pi_k = 1$ . Also note that for  $\{\pi_k\}_{k=1}^K$ , the part of the objective that matters is

$$L(\pi_1, \dots, \pi_K) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{I}[y_n=k] \log \pi_k$$

(Lagrangian form)  
Need to min over  $\pi_i$  =  $\sum_{k=1}^K N_k \log \pi_k$   
max over  $\uparrow$   $\downarrow$  # of points  
with  $y_n=k$

The Constrained formulation will be

$$L(\pi, \lambda) = \underbrace{\sum_{k=1}^K N_k \log \pi_k}_{\text{part of}} + \lambda \left[ \sum_{k=1}^K \pi_k - 1 \right]$$

Taking derivative w.r.t.  $\pi_k$ ,

$$\frac{N_k}{\pi_k} + \lambda = 0$$

$$\pi_k = -\frac{N_k}{\lambda}$$

Also, using  $\sum_{k=1}^K \pi_k = 1$  gives  $\lambda = -N$

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

Now, for  $\mu_k$  and  $\Sigma_k$

The relevant part of the objective function is

$$L(\{\mu_k, \Sigma_k\}) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{I}[y_n=k] \log \mathbf{N}(x_n | \mu_k, \Sigma_k)$$

MLE for  $\mu_k, \Sigma_k$  is exactly the same as the MLE procedure for a multivariate Gaussian, but only using inputs  $x_n$  from class  $k$ .

(August 21 lecture)

I have provided additional slides on the class webpage showing how to do it for multivariate Gaussian. The procedure for this problem will be exactly the same (but only using  $x_n$  with  $y_n=k$ )

## Practice Set 3 (Problem 2)

MAP estimation for  $[\pi_1, \dots, \pi_K]$  is identical  
 ↗ to the MLE procedure, except now we  
 maximize  $\log p(y|\pi) + \log p(\pi)$   
 $\uparrow$   
 (Dirichlet)

$$\log p(y|\pi) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{1}[y_n=k] \log \pi_k$$

$$\log p(\pi) = \log \prod_{k=1}^K \pi_k^{\alpha-1} = \sum_{k=1}^K (\alpha-1) \log \pi_k$$

(Dirichlet without pref. const.)

we also have the sum to 1 constraint.

The full objective (Lagrangian will be)

$$L(\pi, \lambda) = \sum_{k=1}^N N_k \log \pi_k + \sum_{k=1}^K (\alpha-1) \log \pi_k + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

$\downarrow$  log lik       $\downarrow$  log prior  
(wasn't there in Prob. 1)

We will use the same procedure as used in prob. 1.

The solution will be

$$\hat{\pi}_k = \frac{N_k + \alpha-1}{\sum_{k=1}^K (N_k + \alpha-1)}$$

$$\begin{aligned}
 &= \frac{N_k + \alpha-1}{N + K(\alpha-1)} \\
 &= \frac{N_k + \alpha-1}{N + Nd - K}
 \end{aligned}$$

$$\hat{\pi}_k = \frac{N_k + \alpha-1}{N + Nd - K}$$

Alternatively, we can easily get the posterior, which will be Dirichlet due to conjugacy and the form will be  
(multiply the multinoulli and Dirichlet and verify)

$$P(\pi|y) = \text{Dir} [\alpha + N_1, \alpha + N_2, \dots, \alpha + N_k]$$

The MAP estimate will be the mode of this distribution.

$$\hat{\pi}_k = \frac{\alpha + N_k - 1}{N + N\alpha - K}$$

### Practice Set 3 (Problem 3)

Using the expression of Gaussian Conditional from Gaussian joint, we get

$$p(y|x) = N(\mu', \sigma'^2) \text{ where } \mu' \text{ and } \sigma'^2 \text{ are}$$

$$\mu' = \mu_y + V^T \Sigma_{xx}^{-1} (x - \mu_x)$$

$$\sigma'^2 = \sigma_y^2 - V^T \Sigma_{xx}^{-1} V$$

Note that if  $\Sigma_{xx}$  is diagonal then "  $w^T x$ "

$$\mu' = V^T(x - \mu_x) + \mu_y$$

where  $V'$  is  $V$  with each of its entries divided by the corresponding diagonal entry of  $\Sigma_{xx}$ .

Similarity/difference with linear regression:

if  $\mu_y$  and  $\mu_x = 0$  then

$$\mu' = V^T x$$

Probabilistic  
similar to a linear  
regression model's  
mean)

and

$$\sigma'^2 = \sigma_y^2 - V^T \Sigma_{xx}^{-1} V$$

(slightly different  
from prob. linear  
regression's  
variance  
 $P(y|x) = N(w^T x, \beta^2)$ )