

Date \_\_\_\_\_

## ML Assignment - 1

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

Name Mahendra Kumar  
Roll no 1703209

Q. Given, there are only classes (1 &amp; 2)

$$P_1 = P_2 = \frac{1}{2}$$

also, class conditional density's  $f_i$ 's are

$$f_i(x) = \frac{1}{\pi b} \cdot \frac{1}{1 + \left(\frac{x - \alpha_i}{b}\right)^2} \quad i=1,2.$$

(a) To show:  $f_i$ 's are probability density functions.Sol. If  $f_i$ 's are probability density functions then

$$\int_{-\infty}^{\infty} f_i(x) dx = 1.$$

$$\text{LHS} = \int_{-\infty}^{\infty} f_i(x) dx = \int_{-\infty}^{\infty} \frac{1}{\pi b} \cdot \frac{1}{1 + \left(\frac{x - \alpha_i}{b}\right)^2} dx$$

$$= \frac{1}{\pi b} \cdot \left[ \tan^{-1} \left( \frac{x - \alpha_i}{b} \right) \right]_{-\infty}^{\infty}$$

$$= \frac{1}{\pi} \left[ \tan^{-1}(\infty) - \tan^{-1}(-\infty) \right]$$

$$= \frac{1}{\pi} \left( \frac{\pi}{2} - \left( -\frac{\pi}{2} \right) \right) = \frac{\pi}{\pi} = 1. \Rightarrow \text{RHS.}$$

Q.E.D.

$$\text{S. } \int_{-\infty}^{\infty} f_i(x) dx = 1$$

Ans.

(b) Using Bayes rule,

$$\begin{aligned} q_i(x) &= \frac{b_i f_i(x)}{\text{posterior}} \\ &\propto b_i \sum_{i=1}^2 b_i f_i(x) \end{aligned}$$

$$\therefore b_1 = b_2 = b$$

$$\Rightarrow q_i(x) = \frac{\sum_i f_i(x)}{b(f_1(x) + f_2(x))} = \frac{f_i(x)}{f_1(x) + f_2(x)}$$

Decision boundary:

If  $q_1(x) > q_2(x) \Rightarrow x$  belongs to 1<sup>st</sup> class.

$$q_1(x) > q_2(x) \Leftrightarrow \frac{f_1(x)}{f_1(x) + f_2(x)} > \frac{f_2(x)}{f_1(x) + f_2(x)}$$

$$\Rightarrow f_1(x) > f_2(x)$$

$$\Rightarrow \left(\frac{1}{\pi b}\right) \cdot \frac{1}{1 + \left(\frac{x - q_1}{b}\right)^2} > \left(\frac{1}{\pi b}\right) \cdot \frac{1}{1 + \left(\frac{x - q_2}{b}\right)^2}$$

$$\Rightarrow \left(\frac{x - q_1}{b}\right)^2 > \left(\frac{x - q_2}{b}\right)^2$$

$$\Rightarrow x^2 - 2xq_1 + q_1^2 > x^2 - 2xq_2 + q_2^2$$

$$\Rightarrow 2x(q_2 - q_1) \leq q_2^2 - q_1^2$$

$$\textcircled{P} \quad 2x(a_2 - a_1) \leq (a_2 - a_1)(a_2 + a_1).$$

case-1 if  $a_1 = a_2$ . This is always true. i.e. all bomb belongs to class 1.

case-2 if  $a_1 \neq a_2$ . then-

$$2x \leq a_2 + a_1$$

$$\boxed{x \leq \left[ \frac{a_2 + a_1}{2} \right]} \rightarrow \text{decision boundary.}$$

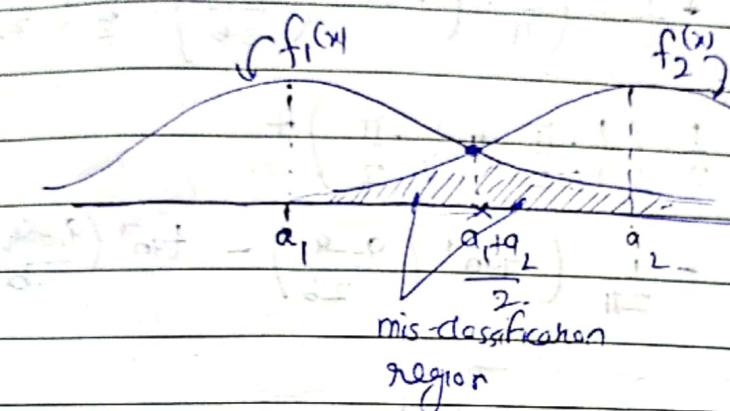
C probability of error for Bayes classifier with 0-1 loss function.

From previous part of it is clear that our bayes classifier is  $h_B(x)$  given by

$$h_B(x) = \begin{cases} 1 & \text{if } x \leq \frac{a_1 + a_2}{2} \\ 2 & \text{if } x > \frac{a_1 + a_2}{2} \end{cases}$$

contd...

(c)



$$\text{Prob} = \text{Prob of misclassification}$$

So probability of error is

$$\text{Error} = P(h_B(x) \neq \text{true class})$$

$$= P(\text{true class} = 1) \cdot P(h_B(x) \neq \text{true class} | \text{true class} = 1)$$

$$+ P(\text{true class} = 2) \cdot P(h_B(x) \neq \text{true class} | \text{true class} = 2)$$

$$= b_1 \cdot P(h_B(x) \neq 1) + b_2 \cdot P(h_B(x) \neq 2)$$

Cars when  $x > \frac{q_1+q_2}{2}$

but  $x$  belongs to class 1.

Cars when  $x \leq \frac{q_1+q_2}{2}$

but  $x$  belongs to class 2.

$$= b_1 \cdot 1 \cdot \int_{\frac{q_1+q_2}{2}}^{\infty} f_1(x) dx + b_2 \cdot 1 \cdot \int_{-\infty}^{\frac{q_1+q_2}{2}} f_2(x) dx + 0 + 0$$

$$= \frac{1}{2} \cdot \left( \frac{b}{\pi b} \left[ \tan^{-1}(\infty) - \tan^{-1}\left(\frac{q_1+q_2 - q_1}{b}\right) \right] \right)$$

$$+ \frac{1}{2} \cdot \left( \frac{b}{\pi b} \left[ \tan^{-1}\left(\frac{q_1+q_2 - q_2}{b}\right) - \tan^{-1}(-\infty) \right] \right)$$

$$\Rightarrow P_{\text{error}} = \frac{1}{2} \cdot \left( \frac{1}{\pi} \left( \frac{\pi}{2} - \tan^{-1} \left( \frac{q_2 - q_1}{Q^b} \right) \right) \right)$$

$$+ \frac{1}{2} \cdot \left( \frac{1}{\pi} \cdot \left( \pi + \tan^{-1} \left( \frac{a_1 - a_2}{2b} \right) - \frac{\pi}{2} \left( -\frac{\pi}{2} \right) \right) \right)$$

$$= \frac{1}{2} \cdot \left( \frac{1}{\pi} \cdot \frac{\pi}{2} + \frac{1}{\pi} \cdot \frac{\pi}{2} \right) +$$

$$-\frac{1}{2\pi} \left( \tan^{-1} \left( \frac{q_2 - q_1}{2b} \right) - \tan^{-1} \left( \frac{q_1 - q_2}{2b} \right) \right)$$

$$\therefore \tan^{-1}(-x) = -\tan^{-1}(x)$$

$$\Rightarrow P_{\text{error}} = \frac{1}{2} \cdot \frac{\pi}{\pi} \left( \frac{\pi}{\pi} \right) - \frac{1}{2\pi} \cdot \left( 2 \tan^{-1} \left( \frac{q_2 - q_1}{2b} \right) \right)$$

$$b_{\text{error}} = \frac{1}{2} \left( \frac{1 + \tan^{-1} \left( \frac{q_2^k - q_1^k}{2b} \right)}{\pi} \right) \quad \underline{\text{Arg.}}$$

Given

Q-2.  $P(h(x)=i|x)$   $\rightarrow$  probability of choosing class i out of a class.

B

① Risk of the classifier on 0-1 loss function

$R(h(x)=i|x) \rightarrow$  risk in choosing class i out of a classes is defined as

$$R(h(x)=i|x) = \sum_{j=1}^a L(h(x)=i, y=j) P(h(x)=j|x)$$

$$= \sum_{j=1}^a p(h(x)=j|x)$$

$$= \sum_{\substack{j=1 \\ j \neq i}}^a p(h(x)=j|x) + 0 \times p(h(x)=i|x)$$

$$R(h(x)=i|x) = 1 - \underbrace{p(h(x)=i|x)}_{(\because \sum_{j=1}^a p(h(x)=j|x) = 1)} P(Y=i|x)$$

②

Current out decision rule is  $\hat{P}$  for given

$\hat{P}(x) \leq 0.5 \rightarrow$  choose class i.s.f.

$$i = \arg \max_j P(h(x)=j|x) \quad \dots$$

$$\Leftrightarrow i = \arg \min_j (-P(h(x)=j|x))$$

$$i = \arg \min_j (1 - P(h(x)=j|x)) \quad \textcircled{O}$$

to get best classifier we want.

risk.  $R(h(x)=i|x)$  to be minimum  
for chosen  $i$  i.e.

$$i = \operatorname{argmin}_j R(h(x)=j|x)$$

$$i = \operatorname{argmin}_j (1 - p(Y=j|x)) \quad \text{--- (2)}$$

Combining (1) & (2) we will get improved classifier if

$$i = \operatorname{argmin}_j (1 - p(Y=j|x)) = \underset{\text{posterior}}{\text{P}(Y=i|x)}$$

$$= \operatorname{argmin}_j (1 - p(h(x)=j|x))$$

$$\text{i.e. } 1 - p(Y=i|x) = 1 - p(h(x)=i|x)$$

$$\Leftrightarrow p(h(x)=i|x) = p(Y=i|x)$$

$\Rightarrow$  i.e. choose class  $i$  out of a class if based on decision rule that

$p(h(x)=i|x)$  is equal to posterior probability of  $Y=i$  given  $x$ .

Q-3 Given, loss function  $L$

$$L(h(x)=i, Y=j) = \begin{cases} 0 & i=j \\ \lambda_r & i=k+1 \rightarrow \text{rejection error} \\ \lambda_m & \text{otherwise} \end{cases}$$

$\lambda_r \rightarrow$  loss for choosing rejection class ( $k+1$ )

$\lambda_m \rightarrow$  loss for mis-classification

$K \rightarrow$  total no. of classes.

10. ~~offset~~

or risk

W $R(h(x)=i|x) \rightarrow$  expected loss associated

with taking action  $i$  for given  $x$  then

Case 1 when  $i=1, 2, \dots, k$  # no. of classes =  $K$

$$R(h(x)=i|x) = \sum_{j=1}^{i-1} L(h(x)=i, Y=j) P(Y=j|x)$$

$$= \sum_{j=1}^{i-1} L(h(x)=i, Y=j) P(Y=j|x) +$$

$$L(h(x)=i, Y=i) P(Y=i|x) +$$

$$\sum_{j=i+1}^K L(h(x)=i, Y=j) P(Y=j|x)$$

$$= \phi \cdot \sum_{\substack{j=1 \\ j \neq i}}^K L(h(x)=i, Y=j) P(Y=j|x)$$

$$= \lambda_m \sum_{\substack{j=1 \\ j \neq i}}^K P(Y=j|x)$$

$$R(h(x)=i|x) = \lambda_m \left( 1 - p(Y=i|x) \right)$$

$$\left( \dots \sum_{j=1}^K p(Y=j|x) = 1 \right)$$

$$\Rightarrow R(h(x)=i|x) = \lambda_m (1 - q_i(x))$$

$\downarrow$  posterior probability  
 $(\propto) \propto P(i)$  for choosing  $i$  given  $x$ .

Case - 2. when  $i = K+1$  i.e. rejection.

$$R(h(x)=K+1|x) = \sum_{j=1}^K L(\underbrace{h(x)=K+1}_{\text{reject}}, Y=j) p(Y=j|x)$$

$$= \lambda_r \sum_{j=1}^K p(Y=j|x)$$

$$R(h(x)=K+1|x) = \lambda_r \dots$$

Now, using bayesian decision theory. we decide a class  $i$  if

$$(i) R(h(x)=i|x) \leq R(h(x)=j|x) \quad \forall j=1, 2, \dots, K$$

$\Rightarrow$  if ~~incorrect~~ i.e. (reject) and.

$$(ii) R(h(x)=i|x) \leq R(h(x)=K+1|x).$$

Simplifying (i) gives

## Simplification - II

Simplification - (1) if  $i < k+1$  then.  
in case-1.

$$\text{case 1. } \lambda_m(1 - q_i(x)) \leq \lambda_m(1 - q_j(x))$$

$$= q_i(x) \geq q_j(x) \quad \text{provided } \lambda_m \neq 0$$

$\cos x - e_2$  if  $i = k+1$

$$\lambda_r \leq \lambda_m (1 - q_j(x))$$

$$\Rightarrow q_j(x) \leq 1 - \frac{\lambda_r}{\lambda_m}$$

6

~~Since~~  $\nabla q_j(x) \geq 0 \Rightarrow$  we choose

~~rejecting~~ rejection class over any other class if &

$$\text{es ist } x^* \text{ ein Optimalpunkt} \quad 1 - \frac{\lambda x}{\lambda_m} \geq q_j(x) \geq 0 \Rightarrow \boxed{\lambda x \leq \lambda_m}$$

## Simplification - ii)

$$R(h(x)=i \mid x) \leq R(h(x)=g \mid x)$$

case-I clearly for  $i=k+1$  equality holds so no need to check.

case II. for  $i < |K+1|$   $\beta_i \in \{\alpha_1, \dots, \alpha_{|K+1|}\}$

$$\text{LHS} \Rightarrow \lambda_m (1 - q_j(\beta_1)) \leq \lambda_Y \cdot 1 \in \text{RHS}$$

$$\Rightarrow q_i(x) \geq 1 - \frac{\Delta}{\lambda_m}$$

Since  $q_i(x) \leq 1$ , this value makes sense if

$$1 - \frac{\lambda_r}{\lambda_m} \leq q_i(x) \leq 1.$$

$$\Rightarrow \left[ \frac{\lambda_r}{\lambda_m} \geq 0 \right] \text{ which is always true.}$$

Part (b) of question.

- What happens if  $\lambda_r = 0$ .

See Simplification - (i) case - 2. If  $\lambda_r = 0$

then this case always hold. i.e. we will always reject. and expected loss will in turn be 0.

- What happens if  $\lambda_r > \lambda_m$ .

See Simplification - (ii) case - II. If  $\frac{\lambda_r}{\lambda_m} > 1$ .

$\Rightarrow q_i(x) \geq 1 - \frac{\lambda_r}{\lambda_m} \Leftrightarrow$  but  $q_i(x) \geq 0$   $\forall i$   
 $\underbrace{\lambda_r}_{\text{non-negative}}$   $\underbrace{\lambda_m}_{\text{non-negative}}$

$\Rightarrow$  this condition holds for every  $i$

So we will definitely not select rejection class.

$\Rightarrow$  So by Simplification - (ii) case - I. we will select class  $i$  having maximum posterior probability i.e.  $q_j(x)$ .

Q-4

## Maximum Likelihood Estimation

## (1) exponential distribution.

the pdf of exponential distribution is

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

Let  $\text{data } \underline{x} = \{x_i\}_{i=1}^N$  be a iid sample from given distribution with unknown parameter  $\lambda$ . Then likelihood will be

$$\text{likelihoold} = \prod f(x_i)$$

$$L(\theta | \lambda) = \prod f(x_1 = x_1, \dots, x_N = x_n)$$

$$L(\theta | \lambda) = \prod_{i=1}^N f(x_i = x_i) \quad (\because x_i \text{ are iid})$$

take log on both sides.

$$\log(L(\theta | \lambda)) = l(\lambda) = \sum_{i=1}^N \log(f(x_i = x_i))$$

$$= \sum_{i=1}^N \log(\lambda e^{-\lambda x_i})$$

$$= \sum_{i=1}^N \log \lambda - \sum_{i=1}^N \lambda x_i$$

$$l(\lambda) = N \log \lambda - \lambda \sum_{i=1}^N x_i \quad \boxed{\text{--- ①}}$$

$$\lambda_{MLE} = \underset{\lambda}{\operatorname{argmax}} \ell(\lambda)$$

taking differentiation on both sides of (1) to  
get w.r.t  $\lambda$  we get -

$$\frac{d(\ell(\lambda))}{d\lambda} = \frac{N}{\lambda} - \sum_{i=1}^N x_i - \varphi$$

equating (2) to zero gives maximum likelihood estimation of  $\lambda$ .

$$\lambda_{MLE} = \frac{1}{N} \left( \sum_{i=1}^N x_i \right).$$

## (2) Multivariate Gaussian Distribution

Assume we have  $n$  random vectors each of size  $p$ :  $x^{(1)}, x^{(2)}, \dots, x^{(n)}$  where each random vector can be interpreted as observation across  $p$  variables. If each  $x^{(i)}$  are i.i.d. or multivariate Gaussian vectors.

$$x^{(i)} \sim \mathcal{N}(\bar{\mu}, \Sigma)$$

$\downarrow$  covariance matrix of shape  $p \times p$

where  $\bar{\mu}, \Sigma$  are unknown parameters.

Notes Notes

Note that by independence of random vector  
 the joint density of data  $\vec{x} = \sum x^{(i)}$   
 is product of individual densities

i.e.  $\prod_{i=1}^N f_{X^{(i)}}(x^{(i)} | \vec{\mu}, \Sigma)$

take logarithm gives log-likelihood function

$$l(\vec{\mu}, \Sigma | \vec{x}) = \log \left( \prod_{i=1}^N f_{X^{(i)}}(x^{(i)} | \vec{\mu}, \Sigma) \right)$$

$$= \log \left( \prod_{i=1}^N \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x^{(i)} - \vec{\mu})^\top \Sigma^{-1} (x^{(i)} - \vec{\mu})} \right)$$

$$= \sum_{i=1}^N \left[ -\frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x^{(i)} - \vec{\mu})^\top \Sigma^{-1} (x^{(i)} - \vec{\mu}) \right]$$

$$l(\vec{\mu}, \Sigma | \vec{x}) = -\frac{N}{2} \log 2\pi - \frac{N}{2} \log |\Sigma|$$

$$- \frac{1}{2} \sum_{i=1}^N (x^{(i)} - \vec{\mu})^\top \Sigma^{-1} (x^{(i)} - \vec{\mu})$$

### Deriving MLE

Take derivative w.r.t.  $\vec{\mu}$  given

$$\frac{\partial}{\partial \vec{\mu}} (l(\vec{\mu}, \Sigma | \vec{x})) = -\frac{1}{2} \left( \sum_{i=1}^N 2\Sigma^{-1} (x^{(i)} - \vec{\mu}) \right)$$

$$\left\{ \therefore \frac{\partial}{\partial \omega} (\omega^T A \omega) = 2 A \omega \text{ if } A \text{ is symmetric} \right\}$$

equating it to 0 gives

$$\begin{aligned} & \sum_{i=1}^N \sum_{j=1}^J (x^{ij} - \bar{x}) = 0 \\ & \sum_{i=1}^N (x^{ij} - \bar{x}) = 0 \quad \left( \text{by taking } \sum_{j=1}^J \text{ on both sides, } \sum_{j=1}^J I = I \rightarrow \text{identity matrix} \right. \\ & \quad \left. \text{and } \sum_{j=1}^J 0 = 0 \rightarrow \text{null matrix} \right). \end{aligned}$$

$$\sum_{i=1}^N x^{ij} - N \bar{x} = 0 \Rightarrow \bar{x} = \frac{1}{N} \sum_{i=1}^N x^{ij}$$

Deriving  $\sum_{MLE}$ .

Deriving MLE for covariance matrix  $\Sigma$ , we will use following results from linear algebra & calculus.

- ① trace is invariant on cyclic permutations of matrix products.

$$\text{tr}[ACB] = \text{tr}[CAB] = \text{tr}[BCA]$$

- ②  $\because x^T A x$  is scalar we take trace
- trace & it remain same.

$$\text{i.e. } \text{tr}(x^T A x) = x^T A x$$

$$\textcircled{3} \quad \frac{\partial}{\partial A} (\text{tr}(AB)) = B^T$$

$$\textcircled{1} \quad \frac{\partial}{\partial A} (\log |A|) = (A^{-1})^T$$

Combining those properties we get -

$$\frac{\partial}{\partial A} (x^T A x) = \frac{\partial}{\partial A} (x^T x A) = (x x^T)^T = x x^T$$

Now back to problem -

We can write log-likelihood function and compute derivative w.r.t  $\Sigma^{-1}$  & using  $|\Sigma^{-1}| = |\Sigma|^{-1}$  we get -

$$\ell(\mu, \Sigma | \mathcal{D}) = -\frac{Np}{2} \log 2\pi + \frac{N}{2} \log |\Sigma^{-1}|$$

$$= \frac{1}{2} \sum_{i=1}^N (x^{(i)} - \bar{x})^T \Sigma^{-1} (x^{(i)} - \bar{x})$$

$$\frac{\partial}{\partial \Sigma} \ell(\mu, \Sigma | \mathcal{D}) = 0$$

$$\Rightarrow 0 + \frac{N}{2} \sum_{i=1}^N -\frac{1}{2} \cdot \sum_{j=1}^N (x_j^{(i)} - \bar{x}_j)(x_j^{(i)} - \bar{x}_j)^T = 0$$

$(\because (\Sigma^{-1})^T = \Sigma^T = \Sigma)$

$$\Rightarrow N \sum_{i=1}^N (x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^T$$

$$\Rightarrow \boxed{\Sigma_{MLE} = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^T}$$

unknown parameter  $\theta := \mu$

Q.5.

Data  $D = \{x_i\}_{i=1}^N$  as iid sampled from normal distribution.  $f_X(x | \mu) \sim N(x; \mu, \sigma^2)$

$$\text{i.e. } f_{X_i}(x_i | \mu) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \quad (I) \quad \text{known}$$

also given  $\mu \sim N(\mu_0, \sigma_0^2)$

$$\text{i.e. } p(\mu) = \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(\mu - \mu_0)^2}{2\sigma_0^2}} \quad (II)$$

this is P.d.f not probability

posterior distribution  $p(\mu | D)$

↑  
this is P.d.f. not probability

(i) Bayesian estimation with parameter corresponding to mode of by bayes rule posterior distribution.

$$p(\mu | D) = \frac{f_X(D | \mu) \cdot p(\mu)}{f_X(D)}$$

∴  $f_X(D)$  is independent of unknown parameter  $\mu$

$$\therefore \hat{\theta}_{MAP} := \underset{\Theta = \mu}{\operatorname{argmax}} p(\mu | D)$$

MAP means maximum a posteriori

mode  $\rightarrow \mu$  at which  $p(\mu | D)$  obtain its maximum value

$$\hat{\mu}_{MAP} = \underset{\Theta = \mu}{\operatorname{argmax}} \frac{f_X(D | \mu) \cdot p(\mu)}{f_X(D)}$$

$$= \underset{\Theta = \mu}{\operatorname{argmax}} f_X(D | \mu) \cdot p(\mu)$$

Assume my samples  $D = \{x_i\}_{i=1}^N$  are ~~sampled~~ iid <sup>are drawn</sup>

$$\hat{\mu}_{MAP} = \underset{\mu}{\operatorname{argmax}} \left( \prod_{i=1}^N f_{x_i}(x_i|\mu) \right) \cdot \log(b(\mu))$$

~~log likelihood~~ <sup>prior</sup>

taking log on both sides, gives.

$$l(\theta) = \underset{\mu}{\operatorname{argmax}} \sum_{i=1}^N \log(f_{x_i}(x_i|\mu)) + \log(b(\mu))$$

$\Rightarrow$  let's call this  $l(\theta)$

$$l(\theta) = \sum_{i=1}^N \log(f_{x_i}(x_i|\mu)) + \log(b(\mu))$$

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

$$+ \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) - \frac{1}{2\sigma^2} (\mu - \mu_0)^2$$

taking differentiation on both sides w.r.t.  $\mu$

and get it to zero gives critical point

$$\frac{\partial l(\mu)}{\partial \mu} = 0 \Rightarrow \left[ \sum_{i=1}^N \frac{1}{\sigma^2} (x_i - \mu) \right] - \frac{1}{\sigma^2} (\mu - \mu_0) = 0$$

$$\Rightarrow \left( \sum_{i=1}^N \frac{1}{\sigma^2} x_i + \frac{\mu_0}{\sigma^2} \right) = \mu \left( \sum_{i=1}^N \frac{1}{\sigma^2} + \frac{1}{\sigma^2} \right)$$

$$\hat{\mu} = \left( \frac{1}{\sigma^2} \sum_{i=1}^N \frac{x_i}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right)$$

$$\hat{\mu}_{MAP} = \frac{1}{N} \left( \sum_{i=1}^N x_i + \frac{\sigma^2 \mu_0}{\sigma_0^2} \right) \quad (11)$$

8. density function using parameters corresponding to mode of posterior distribution is

$$f_X(x|\hat{\mu}_{MAP}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\hat{\mu}_{MAP})^2}{2\sigma^2}}$$

where  $\hat{\mu}_{MAP}$  is as defined in (11)

② Bayesian estimator corresponding to mean of posterior distribution

~~Bayesian Estimator~~

again by choosing - Using ⑦ from previous part

$$p(\mu|\mathcal{D}) = f_X(\mathcal{D}|\mu) \cdot p(\mu)$$

$f_X(\mathcal{D}) \sim$  independent of  $\mu$

$$\propto f_X(\mathcal{D}|\mu) \cdot p(\mu)$$

$$= \text{Constant} \times f_X(\mathcal{D}|\mu) p(\mu)$$

$$p(\mu|\mathcal{D}) = K \cdot f_X(\mathcal{D}|\mu) p(\mu)$$

↓ say

Assuming samples are iid, we get - ..

$$p(\mu | \mathcal{D}) = K \cdot \left( \prod_{i=1}^N p_{x_i}(x_i | \mu) \right) \cdot p(\mu)$$

$$\propto e^{-\frac{(x_1 - \mu)^2}{2\sigma^2}} \cdot e^{-\frac{(x_2 - \mu)^2}{2\sigma^2}} \cdots e^{-\frac{(x_N - \mu)^2}{2\sigma^2}}$$

$$\propto e^{-\frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^2}} = e^{-\frac{(\mu - \bar{x})^2}{2\sigma^2}}$$

proportional  
symbol

$$N \cdot e^{-\frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^2}} = e^{-\frac{\sigma^2 (\mu - \bar{x})^2}{2\sigma^2}}$$

$$\propto e^{-\frac{\sigma^2 (\mu - \bar{x})^2}{2\sigma^2}}$$

$$\text{Simplification we get: } -\mu^2 \left( \sum_{i=1}^N \sigma_0^2 + \sigma^2 \right) + 2\mu \left[ \left( \sum_{i=1}^N x_i \sigma_0^2 \right) + \bar{x} \sigma^2 \right]$$

$$\propto e^{-\frac{\left( \left( \sum_{i=1}^N x_i^2 \sigma_0^2 \right) + \bar{x}^2 \sigma^2 \right)}{2\sigma^2}}$$

$$\propto e^{-\frac{\sigma^2 (\mu - \bar{x})^2}{2\sigma^2}}$$

$$\propto \exp \left\{ -\mu^2 + 2\mu \left[ \frac{\sum x_i \sigma_0^2 + \bar{x} \sigma^2}{\sum \sigma_0^2 + \sigma^2} \right] - \left[ \frac{\sum x_i^2 \sigma_0^2}{\sum \sigma_0^2 + \sigma^2} \right] \right\}$$

Q-  $\bar{x}_1 = \frac{1}{N} \sum_{i=1}^N x_i$

Let  $\mu_1 = \frac{\sum_{i=1}^N x_i \sigma_0^2 + \mu_0 \sigma^2}{\sum_{i=1}^N \sigma_0^2 + \sigma^2}$  &  $\sigma_1^2 = \frac{\sigma_0^2 \sigma^2}{\sum_{i=1}^N \sigma_0^2 + \sigma^2}$

&  $C_1 = \text{constant}_2 = - \left[ \frac{\sum_{i=1}^N x_i \sigma_0^2 + \mu_0 \sigma^2}{\sum_{i=1}^N \sigma_0^2 + \sigma^2} \right]$

then

$$\begin{aligned} p(\mu | D) &\propto \exp \left( \frac{-\mu^2 + 2\mu \mu_1 + C_1}{2\sigma_1^2} \right) \\ &\propto \exp \left( -\frac{(\mu - \mu_1)^2}{2\sigma_1^2} \right) \cdot \exp \left( \frac{C_1 + \mu_1^2}{2\sigma_1^2} \right) \end{aligned}$$

(by adding & subtracting  $\mu_1^2$  in numerator)

$$p(\mu | D) \propto \exp \left( -\frac{(\mu - \mu_1)^2}{2\sigma_1^2} \right)$$

$$p(\mu | D) = \text{constant} \times \exp \left( -\frac{(\mu - \mu_1)^2}{2\sigma_1^2} \right)$$

Clearly  $p(\mu | D)$  is normally distributed s.t  
 $p(\mu | D) \sim N(\mu_1, \sigma_1)$

Let's simplify  $\mu_1$  &  $\sigma_1$  a little

$$\mu_1 = \frac{N}{N} \sum_{i=1}^N x_i \sigma_0^2 + \mu_0 \sigma^2 = \frac{1}{N} \sum_{i=1}^N x_i + \frac{\mu_0 \sigma^2}{\sigma_0^2 N}$$

$$\sum_{i=1}^N \sigma_0^2 + \sigma^2$$

$$\approx 1 + \sigma^2$$

$$\frac{\sigma^2}{N \sigma_0^2}$$

$$\therefore \text{Sample mean } \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\mu_1 = \bar{x} + \frac{\mu_0 \sigma^2 / \sigma_0^2}{1 + (\sigma^2 / \sigma_0^2)}$$

and Variance  $\sigma_1^2$  is

$$\sigma_1^2 = \frac{\sigma_0^2 \sigma_0^2}{N \sigma_0^2 + \sigma^2} = \frac{\sigma_0^2 (\sigma / \sigma_0)^2}{N + (\sigma / \sigma_0)^2}$$

$$\sigma_1^2 = \frac{\sigma^2}{N + (\sigma / \sigma_0)^2}$$

I guess you noticed we have found our answer i.e.

the mean of posterior distribution is

$$\mu_{\text{mean}} = \mu_1 = \frac{\bar{x} + \mu_0 \sigma^2 / \sigma_s^2}{N}$$

for

$$= \frac{1 + \sigma^2 / \sigma_s^2}{N}$$

Key observation:

a)  $N \rightarrow \infty$ .

$$\mu_{\text{mean}} \rightarrow \bar{x} = \mu_{\text{MLE}}$$

$\uparrow$  maximum likelihood estimator.

So density function using mean of posterior distribution is

$$f_X(x | \mu_{\text{mean}}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x - \mu_{\text{mean}})^2}{2\sigma^2}}$$

where  $\mu_{\text{mean}}$  is as defined above.

## Q6. Derive E & M steps for

### (a) Mixture of Gaussians

Let  $x_i \sim N(\mu_j, \sigma^2)$ . In this scenario, our  $x_i$  comes from mixture model with (say)  $K$  mixture components. So we have,

Conditional distribution

$$(X_i | z_i = k) \sim N(\mu_k, \sigma^2)$$

where  $z_i \in \{1, 2, \dots, K\}$   
 ↳ latent variable

So, marginal distribution of  $X_i$  is

$$P(X_i = x_i) = \sum_{j=1}^K P(X_i = x_i | z_i = k_j) P(z_i = j)$$

$$P(X_i = x_i) = \sum_{j=1}^K \pi_{kj} P(X_i = x_i | z_i = j)$$

Similarly, Joint probability of observation  $x_1, x_2, \dots, x_N$  assuming iid sample is

$$P(X_1 = x_1, \dots, X_N = x_N) = \prod_{i=1}^N \left( \sum_{j=1}^K \pi_{kj} P(X_i = x_i | z_i = j) \right)$$

$$\Rightarrow f_{x_1, x_2, \dots, x_N}(x_1, x_2, \dots, x_N) = \prod_{i=1}^N \left( \sum_{j=1}^K \pi_{kj} f_{X_i}(x_i | \mu_j, \sigma_j^2) \right) \quad (1)$$

$$\text{where } f_{X_i}(x_i | \mu_j, \sigma_j^2) = \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(x_i - \mu_j)^2}{2\sigma_j^2}}$$

Aim: find maximum likelihood estimates  $\pi_j, \mu_j$  &  $\sigma_j^2$   
 given a dataset  $\sum_{i=1}^N x_i = x_i, i \in N$

Optimization problem

Likelihood function  $L(\theta)$  is

$$L(\theta) = \prod_{i=1}^N \sum_{j=1}^K \pi_j f_{X_i}(x_i | \mu_j, \sigma_j^2) = \prod_{i=1}^N \sum_{j=1}^K \pi_j f_j$$

for notational convenience.

Log Likelihood will be

$$\ell(\theta) = \sum_{i=1}^N \log \left( \sum_{j=1}^K \pi_j f_j \right)$$

$$\text{MLF} \quad \hat{\theta}_{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \ell(\theta) \quad \text{s.t.} \quad \sum_{j=1}^K \pi_j = 1$$

$$L(\theta, \lambda) = \sum_{i=1}^N \log \left( \sum_{j=1}^K \pi_j f_j \right) + \lambda \left( 1 - \sum_{j=1}^K \pi_j \right)$$

$$\frac{\partial L(\theta, \lambda)}{\partial \pi_j} = 0 \Rightarrow \left[ \sum_{i=1}^N \frac{f_j}{\sum_{j=1}^K \pi_j f_j} - \lambda \right] = 0 \quad \text{--- (1)}$$

$$\frac{\partial L(\theta, \lambda)}{\partial \mu_j} = 0 \Rightarrow \left[ \sum_{i=1}^N \frac{\pi_j f_j}{\sum_{j=1}^K \pi_j f_j} \cdot \left( \frac{-1}{2} \cdot \frac{2(x_i - \mu_j)}{\sigma_j^2} \cdot (+) \right) = 0 \right]$$

$$\frac{\partial L(\theta, \lambda)}{\partial \sigma_j^2} = 0 \Rightarrow \left[ \sum_{i=1}^N \frac{\pi_j f_j}{\sum_{j=1}^K \pi_j f_j} \cdot \left( -\frac{1}{\sigma_j^2} f_j + f_j \cdot \left( \frac{-2}{\sigma_j^2} \right) \left( \frac{-1}{2} (x_i - \mu_j)^2 \right) \right) = 0 \right] \quad \text{--- (2)}$$

let's denote:

$$\frac{\pi_j f_j}{\sum_{k=1}^K \pi_k f_k} = \gamma_{ij}$$

then by ①, ② & ③ we have.

$$② \Rightarrow \sum_{i=1}^N \gamma_{ij} (x_i - \mu_j) = 0 \Rightarrow \mu_j = \frac{\sum_{i=1}^N \gamma_{ij} x_i}{\sum_{i=1}^N \gamma_{ij}}$$

$$③ \Rightarrow \sum_{i=1}^N \gamma_{ij} \left( \frac{(x_i - \mu_j)^2}{\sigma_j^2} - 1 \right) = 0 \Rightarrow \sigma_j^2 = \frac{\sum_{i=1}^N \gamma_{ij} (x_i - \mu_j)^2}{\sum_{i=1}^N \gamma_{ij}}$$

$$① \Rightarrow \sum_{i=1}^N \frac{\gamma_{ij}}{\pi_j} - \lambda = 0 \Rightarrow \sum_{i=1}^N \gamma_{ij} - \lambda \pi_j = 0$$

taking sum over j. on both sides given.

$$\sum_{j=1}^K \sum_{i=1}^N \frac{\gamma_{ij}}{\pi_j} - \lambda \sum_{j=1}^K \pi_j = 0 \Rightarrow \sum_{i=1}^N \left( \sum_{j=1}^K \frac{\pi_j f_j}{\pi_1 f_1 + \dots + \pi_K f_K} \right) - \lambda \cdot 1 = 0$$

( $\because \sum_{j=1}^K \pi_j = 1$ )

$$\Rightarrow \sum_{i=1}^N 1 - \lambda = 0 \Rightarrow \boxed{\lambda = N}$$

$$\Rightarrow \pi_j = \sum_{i=1}^N \gamma_{ij} \rightarrow$$

$$\boxed{\pi_j = \frac{\sum_{i=1}^N \gamma_{ij}}{N}}$$

Init: Randomly choose  $\Theta^{(0)} = (\pi_j, \mu_j^{(0)}, \sigma_j^{(0)}) \forall j=1, 2, \dots, K$ .

E-step: -

$$\gamma_{ij}^{(t+1)} = \frac{\pi_j f_j(x_i | \Theta^{(t)})}{\pi_1 f_1 + \dots + \pi_K f_K}$$

M-step

$$\pi_j^{(t+1)} = \frac{\sum_{i=1}^N \gamma_{ij}^{(t+1)}}{N}$$

$$\mu_j^{(t+1)} = \frac{\sum_{i=1}^N \gamma_{ij}^{(t+1)} x_i}{\sum_{i=1}^N \gamma_{ij}^{(t+1)}}$$

$$\sigma_j^{(t+1)} = \sqrt{\frac{\sum_{i=1}^N \gamma_{ij}^{(t+1)} (x_i - \mu_j^{(t+1)})^2}{\sum_{i=1}^N \gamma_{ij}^{(t+1)}}}$$

loop till convergence.

## (b) mixture of Bernoulli

Let  $X_i \sim \text{Bernoulli}(p)$ . In this scenario our  $X_i$  can come from mixture model with (say)  $K$  mixture components so we have conditional distribution

$$(X_i | Z_i = k) \sim \text{Bernoulli}(p_k)$$

where  $Z_i \in \{1, 2, \dots, K\}$

↳ latent variable

so marginal distribution of  $X_i$  is

$$\begin{aligned} P(X_i = x_i) &= \sum_{j=1}^K P(X_i = x_i | Z_i = j) \cdot P(Z_i = j) \\ &= \sum_{j=1}^K \pi_j \cdot P(X_i = x_i | Z_i = j) \end{aligned}$$

Similarly joint probability distribution of observation  $X_1, X_2, \dots, X_N$  assuming iid sample is

$$P(X_1 = x_1, \dots, X_N = x_N) = \prod_{i=1}^N \left( \sum_{j=1}^K \pi_j P(X_i = x_i | Z_i = j) \right)$$

for Bernoulli distribution.

$$P(X_i = x_i | Z_i = j) = p_j^{x_i} (1-p_j)^{1-x_i}$$

Aim: Find maximum likelihood estimator  $\log \pi_j$ 's,  $b_j$ 's

given data  $x_i = x_i^i$ ,  $s_{i=1}^N$

likelihood function  $L(\theta)$  is

$$L(\theta) = \prod_{i=1}^N \left( \sum_{j=1}^K \pi_j b_j^{x_i} (1-b_j)^{1-x_i} \right)$$

log likelihood will be:

$$\ell(\theta) = \sum_{i=1}^N \log \left( \sum_{j=1}^K \pi_j b_j^{x_i} (1-b_j)^{1-x_i} \right)$$

$$\theta_{MLE} = \underset{\pi_j, b_j}{\operatorname{argmax}} \ell(\theta) \quad \text{s.t. } \sum_{j=1}^K \pi_j = 1$$

$$L(\theta, \lambda) = \sum_{i=1}^N \log \left( \sum_{j=1}^K \pi_j b_j^{x_i} (1-b_j)^{1-x_i} \right) + \lambda \left( 1 - \sum_{j=1}^K \pi_j \right)$$

$$\lambda \left( 1 - \sum_{j=1}^K \pi_j \right) = 0$$

$$\frac{\partial L(\theta, \lambda)}{\partial \pi_j} = \sum_{i=1}^N \frac{\pi_j f_j}{\pi_1 f_1 + \pi_2 f_2 + \dots + \pi_K f_K} - \lambda = 0 \quad \boxed{①}$$

$$\text{where } f_j := b_j^{x_i} (1-b_j)^{1-x_i}$$

$$\begin{aligned} \frac{\partial L(\theta, \lambda)}{\partial b_j} &= \sum_{i=1}^N \frac{\pi_j}{\pi_1 f_1 + \dots + \pi_K f_K} \cdot \left( x_i \pi_j b_j^{x_i-1} (1-b_j)^{1-x_i} + (1-x_i) \pi_j b_j^{x_i} (1-b_j)^{1-x_i} \right) \\ &= 0 \end{aligned}$$

$$\sum_{i=1}^N \frac{\pi_i b_j x_i (1-b_j)^{1-x_i}}{\pi_1 f_1 + \dots + \pi_K f_K} \left\{ x_i b_j^{-1} - (1-x_i)(1-b_j)^{-1} \right\} = 0$$

$$\sum_{i=1}^N \frac{\pi_i f_i}{\pi_1 f_1 + \dots + \pi_K f_K} \left\{ x_i \frac{(1-b_j)}{b_j} - (1-x_i) \right\} = 0$$

$$\text{Let } \boxed{\gamma_{ij} = \frac{\pi_j f_i}{\pi_1 f_1 + \dots + \pi_K f_K}} \text{ then}$$

$$\frac{\partial L(\theta, \lambda)}{\partial b_j} = \sum_{i=1}^N \gamma_{ij} \left\{ \frac{x_i}{b_j} - x_i' - 1 + x_i \right\} = 0$$

$$\hat{b}_j = \frac{\sum_{i=1}^N \gamma_{ij} x_i}{\sum_{i=1}^N \gamma_{ij}} \quad \text{--- (Q)}$$

from (Q) we also have

$$\sum_{i=1}^N \frac{f_i}{\pi_1 f_1 + \dots + \pi_K f_K} - \lambda = 0 \Rightarrow \sum_{i=1}^N \frac{\pi_i Y_{ij}}{\pi_j} - \lambda = 0$$

$$\Rightarrow \sum_{i=1}^N \gamma_{ij} - \lambda \pi_j = 0 \quad \begin{matrix} \text{taking sum over } j \text{ on} \\ \text{both sides gives} \end{matrix}$$

$$\sum_{j=1}^K \sum_{i=1}^N \gamma_{ij} - \lambda \sum_{j=1}^K \pi_j = 0 \Rightarrow \lambda = \sum_{i=1}^N \sum_{j=1}^K \frac{\pi_j f_i}{\pi_j f_i + \text{other terms}}$$

$$\Rightarrow \boxed{\lambda = N}$$

$$\text{so } \hat{\pi}_j = \sum_{i=1}^N r_{ij} \Rightarrow$$

$$\hat{\pi}_j = \frac{\sum_{i=1}^N r_{ij}}{N}$$

## EM Algorithm

Init: Randomly choose  $\Theta^0 = (\pi_j^{(0)}, b_j^{(0)})$  for  $j=1, 2, \dots, K$

$$\text{E-step } r_{ij}^{(t+1)} = \frac{\pi_j^{(t)} p_j^{x_i} (1-p_j)^{1-x_i}}{\sum_{k=1}^K \pi_k^{(t)} p_k^{x_i} (1-p_k)^{1-x_i}}$$

 $t+1$ 

$$\text{M-step } \hat{\pi}_j = \frac{\sum_{i=1}^N r_{ij}^{(t+1)}}{N}$$

$$\hat{b}_j^{(t+1)} = \frac{\sum_{i=1}^N r_{ij}^{(t+1)} x_i}{\sum_{i=1}^N r_{ij}^{(t+1)}}$$

loop till convergence

Q-7. Bernoulli distribution, let,  $D = \{x_i\}_{i=1}^N$

$$f(\{x_i\} | b) = b^{x_i} (1-b)^{1-x_i}$$

given conjugate prior to be Beta distribution  
i.e.

$$\begin{aligned} g(b) &\sim \text{Beta}(\alpha, \beta) \\ \Leftrightarrow g(b) &= \text{constant} \propto b^{\alpha-1} (1-b)^{\beta-1} \end{aligned}$$

by bayes rule

$$f(p | D) = \frac{f_X(D | p) \cdot g(p)}{f_X(D)}$$

$f_X(D)$  is independent of  $p$ . and is just a scaling constant.

MAP estimate of Bernoulli distribution i.e. Bernoulli ( $p$ )

$$\begin{aligned} p_{MAP} &= \underset{p}{\operatorname{argmax}} f(p | D) \\ &= \underset{p}{\operatorname{argmax}} \left( \prod_{i=1}^N p^{x_i} (1-p)^{1-x_i} \right) \cdot g(p) \end{aligned}$$

taking log won't change argmax so,

$$p_{MAP} = \underset{p}{\operatorname{argmax}} \sum_{i=1}^N \log (f(x_i | p)) + \log g(p)$$

$$\text{Let } l(b) = \sum_{i=1}^N \log(f(x_i|b)) + \log(g(b))$$

$$= \sum_{i=1}^N \{x_i \log b + (1-x_i) \log(1-b)\} +$$

$$\log(\text{constant}) + (\alpha-1) \log b + (\beta-1) \log(1-b)$$

By bayes rule.

take differentiation w.r.t.  $b$  on both sides

$$\frac{\partial l(b)}{\partial b} = 0 \Rightarrow \sum_{i=1}^N \left( \frac{x_i}{b} - \frac{(1-x_i)}{1-b} \right) + \frac{\alpha-1}{b} - \frac{(\beta-1)}{1-b} = 0$$

"on simplification we get"

$$\boxed{\hat{b}_{\text{MAP}} = \frac{N + (\alpha-1) + (\beta-1)}{\sum_{i=1}^N x_i + (\alpha-1)}} \text{ or } \boxed{\hat{b}_{\text{MAP}} = \frac{\sum_{i=1}^N x_i + (\alpha-1)}{N + (\alpha-1) + (\beta-1)}}$$

Ans

Bayesian estimate of

Q-8

## Multinomial distribution

multinomial distribution.

$$\text{Defn: } \{x_i\}_{i=1}^N$$

$$S: X = (x_1, x_2, \dots, x_K)$$

$$\text{see } \sum_{j=1}^K x_{ij} = m$$

$$D = \{Y_i\}_{i=1}^N$$

$$\text{where } Y_i = \{x_{ij}\}_{j=1}^K$$

$$\sum_{i=1}^K x_{ii} = m \text{ (say)}$$

then probability mass function of multinomial distribution is

$$f(x_1^{(i)}, x_2^{(i)}, \dots, x_K^{(i)} | p_1, p_2, \dots, p_K) = P(X_1^{(i)}=x_1^{(i)}, X_2^{(i)}=x_2^{(i)}, \dots, X_K^{(i)}=x_K^{(i)})$$

this sample

$$= \frac{n!}{x_1^{(i)}! x_2^{(i)}! \dots x_K^{(i)}!} \cdot p_1^{x_1^{(i)}} p_2^{x_2^{(i)}} \dots p_K^{x_K^{(i)}}$$

$$= \text{Constant} \times p_1^{x_1^{(i)}} p_2^{x_2^{(i)}} \dots p_K^{x_K^{(i)}}$$

$$\text{S.t. } \sum x_i = m, x_i \geq 0$$

Given conjugate prior is dirichlet distribution i.e

$$g(p_1, p_2, \dots, p_K) \sim \text{dirichlet}(\alpha_1, \alpha_2, \dots, \alpha_K)$$

$$\Leftrightarrow g(p_1, p_2, \dots, p_K) = \text{(constant)} \times \prod_{i=1}^K p_i^{\alpha_i - 1}$$

$$\text{S.t. } \sum_{i=1}^K p_i = 1 \quad \& \quad p_i \geq 0$$

$$\& i = 1, 2, \dots, K$$

MAP estimate of multinomial distribution

$$\text{Def. } \vec{b} = (b_1, b_2, \dots, b_K) \text{ s.t. } \sum_{i=1}^K b_i = 1$$

$$\vec{b}_{\text{MAP}} = \underset{\vec{p}}{\operatorname{argmax}} f(\vec{p} | \vec{x}) \text{ s.t. } \sum_{i=1}^K p_i = 1$$

by now we know that this is equivalent to

$$\vec{b}_{\text{MAP}} = \underset{\substack{\vec{p} \\ b_1, b_2, \dots, b_K}}{\operatorname{argmax}} \left( \sum_{i=1}^N \log (f(x_i^{(i)}, x_{(i)}^{(i)}, b_1, \dots, b_K)) + \log (b_1, \dots, b_K) \right)$$

Using Lagrange's multiplier

$$\Rightarrow \text{let } l(\lambda b_1, b_2, \dots, b_K) = \sum_{i=1}^N \log (f(x_i^{(i)}, x_{(i)}^{(i)}, b_1, \dots, b_K))$$

$$+ \log (g(b_1, \dots, b_K)) + \lambda \left( 1 - \sum_{j=1}^K b_j \right)$$

$$= \sum_{i=1}^N \log (c_{ij}) + \sum_{i=1}^N \sum_{j=1}^{K(i)} x_j^{(i)} \log b_j + \log (\text{constant}) + \sum_{j=1}^K \log (x_j - 1) \log b_j + \lambda \left( 1 - \sum_{j=1}^K b_j \right)$$

Now

$$\frac{\partial l}{\partial b_j} = 0 \Rightarrow \sum_{i=1}^N \frac{x_j^{(i)}}{b_j} + \frac{(x_j - 1)}{b_j} - \lambda = 0$$

$$\Rightarrow b_j = \frac{\sum_{i=1}^N x_j^{(i)} + (x_j - 1)}{\lambda} \quad \text{--- (1)}$$

$$\lambda b_j = \sum_{i=1}^N x_j^{(i)} + (k_j - 1)$$

Take summation on both sides w.r.t  $j$

$$\sum_{j=1}^k \lambda b_j = \sum_{j=1}^k \sum_{i=1}^N x_j^{(i)} + \sum_{j=1}^k (k_j - 1)$$

$$\lambda \left( \sum_{j=1}^k b_j \right) = \sum_{i=1}^N \left( \sum_{j=1}^k x_j^{(i)} \right) + \sum_{j=1}^k (k_j - 1)$$

sum

$$\boxed{\lambda = mN + \sum_{j=1}^k (d_j - 1)} \quad \text{--- (2)}$$

Substitute (2) in (1) we get

$$\hat{b}_{\text{MAP}} = \frac{\sum_{i=1}^N x_j^{(i)} + (d_j - 1)}{mN + \sum_{j=1}^k (d_j - 1)}$$

$\forall j = 1, 2, \dots, k$

$$\text{So } \hat{b}_{\text{MAP}} = (\hat{b}_1, \hat{b}_2, \dots, \hat{b}_k) \quad \underline{\text{ans}}$$

Required MAP estimation of multinomial distribution

Q-1

## generalized EM for MAP estimation

Sol:

The expectation maximization algorithm, or EM algorithm is a general technique for finding maximum likelihood solutions for probabilistic models having latent variables.

Let  $X \rightarrow$  set of all observed variables.

$Z \rightarrow$  set of all hidden/latent variables.

$\theta \rightarrow$  set of all <sup>unknown</sup> parameters

$D = \{X_i\}_{i=1}^n, q \sim \text{distribution of } Z$

$$\theta_{MAP} = \operatorname{argmax}_{\theta} p(\theta | D) \quad \text{using Bayes rule.}$$

$$= \operatorname{argmax}_{\theta} \frac{p(\theta, D)}{p(D)} \cdot \log \frac{p(\theta, D)}{p(D)} \quad \text{posterior} \quad \text{prior distribution}$$

$$= \operatorname{argmax}_{\theta} \log \left( \frac{p(\theta, D)}{p(D)} \right)$$

( $\because \log$  is increasing function)

$$= \operatorname{argmax}_{\theta} \log p(\theta, D) - \operatorname{argmax}_{\theta} \log p(D) \quad \text{independent of } \theta$$

$$= \operatorname{argmax}_{\theta} \log p(\theta, D)$$

$$= \operatorname{argmax}_{\theta} \log [p(D|\theta) \cdot p(\theta)] \quad \text{prior} \quad \text{likelihood}$$

$$= \operatorname{argmax}_{\theta} (\log p(D|\theta) + \log p(\theta))$$

assuming iid samples.

$$p(D|\theta) = \prod_{i=1}^N p(x_i|\theta)$$

$$\Rightarrow \theta_{MAP} = \underset{\theta}{\operatorname{argmax}} : \sum_{i=1}^N \log p(x_i|\theta) + \log(p(\theta))$$

$$\text{or } L(\theta) = \sum_{i=1}^N \log(p(x_i|\theta)) + \log(p(\theta))$$

$$= \sum_{i=1}^N \log \left( \sum_c p(x_i, z_i=c | \theta) \right) + \log(p(\theta))$$

(~~assume~~ note  $z_i$  is hidden variable.)  
 $\sum_c p(x_i, z_i=c | \theta) = p(x_i | \theta)$

$$= \sum_{i=1}^N \log \left( \sum_c q(z_i=c) \underbrace{\frac{p(x_i, z_i=c | \theta)}{q(z_i=c)}}_{} \right) + \log(p(\theta))$$

(where  $\sum_c q(z_i=c)$  is a probability)

g.e.

$$\sum_c q(z_i=c) = 1, q(z_i=c) > 0$$

$$\geq \sum_{i=1}^N \sum_c q(z_i=c) \log \left( \frac{p(x_i, z_i=c | \theta)}{q(z_i=c)} \right)$$

$$+ \log(p(\theta))$$

(Using Jensen's inequality: i.e.  
 $\log(E(\cdot)) \geq E(\log(\cdot))$ )

$\Rightarrow l(\theta) \geq l(\theta, q)$  decreasing (why?)  $\forall q$  of  $q(z_i=c) > 0$

$\downarrow$  just a notation

$\& \sum_c q(z_i=c) = 1$

E-step: find  $q^*$  s.t.

$$l(\theta) = l(\theta, q^*) - \log(p(\theta))$$

$\Rightarrow$  for any  $q$ , since

$$\therefore l(\theta) - (l(\theta, q)) = \left( \sum_{i=1}^N p \log (p(x_i|\theta) + \log(p(\theta)) \right)$$

$$- \left( \sum_{i=1}^N \sum_c q(z_i=c) \log \left( \frac{p(x_i, z_i=c|\theta)}{q(z_i=c)} \right) \right)$$

$$(l(\theta)) = \log(p(\theta))$$

$$= \sum_{i=1}^N \left( \sum_c q(z_i=c) \log \left( \frac{p(x_i|\theta)}{q(z_i=c)} \right) - \sum_{i=1}^N \sum_c q(z_i=c) x_i \log \left( \frac{p(x_i, z_i=c)}{q(z_i=c)} \right) \right)$$

$$= \sum_{i=1}^N \sum_c q(z_i=c) \log \left( \frac{p(x_i|\theta) q(z_i=c)}{p(x_i, z_i=c|\theta)} \right)$$

$$= \sum_{i=1}^N \left[ \sum_c q(z_i=c) \log \left( \frac{q(z_i=c)}{p(z_i=c|x_i, \theta)} \right) \right]$$

$\hookrightarrow$  KL-divergence

b/w  $q(z_i|c)$  &  $p(z_i|x_i, \theta)$

$$= \sum_{i=1}^N KL(q(z_i) || p(z_i|x_i, \theta))$$

for for E-step find  $q^*$  s.t.

$\text{So } \ell(\theta) = \ell(\theta, q^*) \geq 0 ; \text{ where}$

$$\Rightarrow \sum_{i=1}^N \text{KL}(q^*(z_i) || P(z_i | x_i, \theta)) = 0$$

but we know  $\text{KL}(p || q) \geq 0 \quad \forall p, q$

&  $\text{KL}(p || q) = 0 \iff p = q$

$$\Rightarrow q^*(z_i) := p(z_i | x_i, \theta) \quad \begin{array}{l} \text{independent} \\ \text{of prior } b(\theta) \end{array}$$

E-step  $q^*(z_i) = p(z_i | x_i, \theta)$

M-step find  $\theta^* = \arg \max_{\theta} \ell(\theta, q^*)$

$$\theta^* = \arg \max_{\theta} \left\{ \sum_{i=1}^N \sum_c q(z_i=c) \ln \left( \frac{P(x_i, z_i=c | \theta)}{q(z_i=c)} \right) + \log(b(\theta)) \right\}$$

$$\theta^* = \arg \max_{\theta} \left\{ \sum_{i=1}^N \sum_c q(z_i=c) \ln \left( \frac{P(x_i, z_i=c | \theta)}{q(z_i=c)} \right) + \ln(b(\theta)) \right\}$$

$\because q(z_i=c)$  is independent of  $\theta$

Note  
 $\theta^*$  depends on prior  $b(\theta)$

\* Proof of correctness -

$$l(\theta^{(old)}) \geq \sum_{i=1}^N \sum_{c=1}^C q^*(z_i=c) \ln$$

$$\therefore l(\theta^{(old)}) = l(\theta^{(old)}, q^*) + \sum_{i=1}^N \text{KL}(q^*(z_i) || P(z_i | x_i; \theta))$$

$$( \because q^*(z_i) := p(z_i | x_i; \theta) )$$

$$\therefore l(\theta^{(old)}) = l(\theta^{(old)}, q^*)$$

$$= \sum_{i=1}^N \sum_{c=1}^C q^*(z_i=c) \log \left( \frac{P(x_i, z_i=c)}{q^*(z_i=c)} \right) + \log(b(\theta))$$

$$l(\theta^{(old)}) = \sum_{i=1}^N \sum_{c=1}^C q^*(z_i=c) \log(P(x_i, z_i=c)) + \log(b(\theta))$$

$$l(\theta^{(old)}) = \sum_{i=1}^N \sum_{c=1}^C \log(q^*(z_i=c))$$

$$l(\theta^{(new)}) = l(\theta^{(new)}, q^*) + \sum_{i=1}^N \text{KL}(q^*(z_i) || P(z_i | x_i; \theta^{(new)}))$$

$$= \sum_{i=1}^N \sum_{c=1}^C q^*(z_i=c) \ln(P(x_i, z_i=c)) - \sum_{i=1}^N \sum_{c=1}^C q^*(z_i=c) \log(P(z_i | x_i, \theta^{(new)}))$$

$$+ \log(b(\theta)) + \text{KL}(q^*(z_i) || P(z_i | x_i, \theta^{(new)}))$$

$$= l(\theta^{(new)}) - l(\theta^{(old)}) = \text{KL}(q^*(z_i) || P(z_i | x_i, \theta^{(new)}))$$

$$\therefore l(\theta^{(new)}) \geq l(\theta^{(old)}) \Rightarrow \text{H.P.}$$

END OF ASSIGNMENT

X X X