

1/2/22 ML Assignment - 1

classmate

Date _____
Page _____

Name Mahendra Kumar
Roll no 1703209

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

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

also, class conditional density's f_i 's are

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

(a) To show: f_i 's are probability density functions.

Sol. If f_i 's are probability density function then

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

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

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

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

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

Q.E.D.

$$\text{Q} \quad \int_{-\infty}^{\infty} f_i(x) dx = 1 \quad \underline{\underline{\text{Ans}}}$$

(b) Using Bayes rule,

$$\text{posterior } q_i(x) = \frac{b_i f_i(x)}{\sum_{i=1}^2 b_i f_i(x)}$$

$$\therefore b_1 = b_2 = \frac{1}{2}$$

$$\Rightarrow q_i(x) = \frac{\frac{1}{2} f_i(x)}{\frac{1}{2}(f_1(x) + f_2(x))} = \frac{f_i(x)}{f_1(x) + f_2(x)}$$

Decision boundary:

If $q_1(x) \geq q_2(x) \Rightarrow x$ belongs to 1st class.

$$q_1(x) \geq q_2(x) \Leftrightarrow \frac{f_1(x)}{f_1(x) + f_2(x)} \geq \frac{f_2(x)}{f_1(x) + f_2(x)} \quad \lambda \neq 0$$

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

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

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

$$\text{Therefore: } x^2 - 2xq_2 + q_2^2 \geq x^2 - 2xq_1 + q_1^2$$

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

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

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

case-2 if $a_1 \neq a_2$ then

$$2x \leq a_2 + a_1$$

$$x \leq \left(\frac{a_2 + a_1}{2} \right) \rightarrow \text{decision boundary.}$$

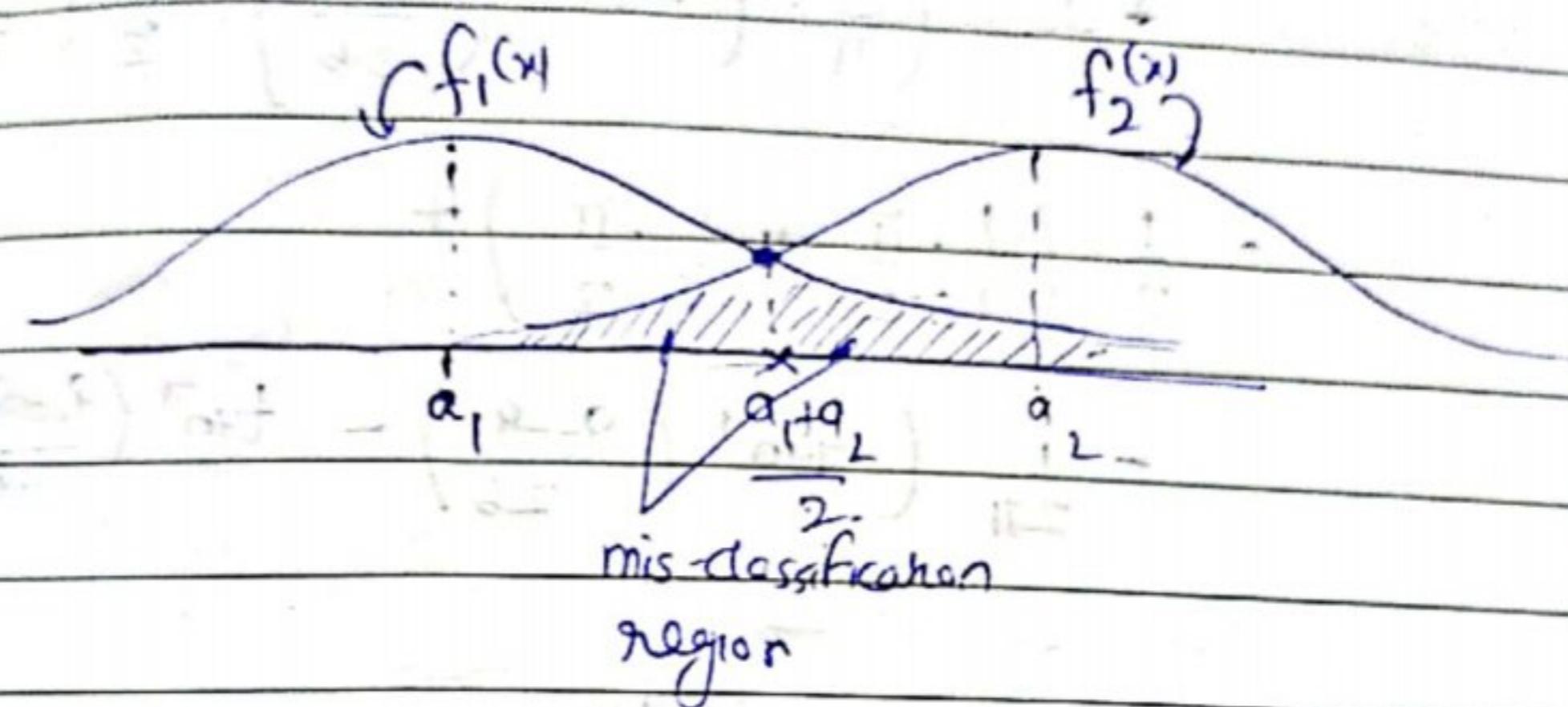
(c) probability of error for Bayes classifier with 0-1 loss function.

From previous part 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)



$$\text{Perror} = \underbrace{\int_{q_1}^{q_2} f(x) dx}_{\text{area under } f(x)} \text{ Perror}$$

So probability of error is

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

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

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

$$= b_1 \cdot \underbrace{P(h_B(x) \neq 1)}_{\text{case when } x > \frac{q_1+q_2}{2}} + b_2 \cdot \underbrace{P(h_B(x) \neq 2)}_{\text{case when } x \leq \frac{q_1+q_2}{2}}$$

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

but x belongs to class 1.

but x belongs to class 2.

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

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

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

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

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

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

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

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

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

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

Given

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

B

@ Risk of this classifier on 0-1 loss function

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

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

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

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

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

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



Current out decision rule is \bullet for given

x choose i s.f.

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

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

$$i = \arg \min_j (1 - P(h(x)=j|x))$$

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

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

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

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

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

\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 $Y=j$ given x .

Q-3 Given, loss function L

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

$\lambda_x \rightarrow$ loss for choosing rejection class $(K+1)$

$\lambda_m \rightarrow$ loss for mis-classification

$K \rightarrow$ total no. of classes.

R or risk

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

with taking action i for given x then.

Case-1 When $i=1, 2, \dots, K$ #no of classes = K

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

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

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

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

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

$\lambda_m + j \& i=1, 2, \dots, K$.

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

$$R(h(x)=i|x) = \lambda_m (1 - P(Y=i|x))$$

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

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

(i.e.) posterior probability
for choosing i given x .

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

$$R(h(x)=K+1|x) = \sum_{j=1}^K L(\underbrace{h(x)=K+1, Y=j}_{\Downarrow}) P(Y=j|x) / \lambda_{K+1}$$

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

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

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 $i < K+1$ i.e. if i and.

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

Simplifying (ii) gives

Simplification -ii)

if $i < k+1$ then.

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

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

case -2. if $i = k+1$.

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

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

but

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

rejection class over any other class if &

$$1 - \frac{\lambda_r}{\lambda_m} \geq q_j(x) \geq 0 \Rightarrow \boxed{\lambda_r \leq \lambda_m}$$

Simplification -iii)

$$R(h(x)=i | x) \leq R(h(x)=k+1 | x)$$

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

$$\text{LHS} \Rightarrow \lambda_m (1 - q_i(x)) \leq \lambda_r \cdot 1. \subset \text{RHS}$$

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

Since $q_i(x) \leq 1$ thus value make sense if

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

$$\Rightarrow \left[\frac{\lambda_r}{\lambda_m} \geq 0 \right] \quad \text{where } \lambda_m \neq 0.$$

Part (b) of question.

- What happen if $\lambda_r = 0$.

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

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

- What happen if $\lambda_r > \lambda_m$

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

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

∴ this condition holds for every i

So we will definitely not select rejection class.

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

Q-4Maximum Likelihood Estimation(1) exponential distribution

the pdf of exponential distribution is

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x > 0 \\ 0 & x \leq 0. \end{cases}$$

Let take $\mathcal{D} = \{x_i\}_{i=1}^N$ be a iid sample from given distribution with unknown parameter λ . Then likelihood will be

$$L(\mathcal{D} | \lambda) = \prod_{i=1}^N f(x_i; \lambda)$$

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

take log on both sides.

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

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

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

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

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

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

$$1) \frac{d(\ell(\lambda))}{d\lambda} = \frac{N}{\lambda} - \sum_{i=1}^N x_i - 0$$

equating $\textcircled{1}$ to zero gives maximum likelihood estimation of λ .

$$\lambda_{MLE} = \frac{N}{\sum_{i=1}^N x_i}$$

$$\lambda_{MLE} = \frac{N}{\sum_{i=1}^N x_i}$$

③ 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.

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

↳ covariance matrix of shape $p \times p$

↳ here \bar{u}, Σ are unknown parameters.

Note Note

Note that by independence of random vector
 the joint density of data $\mathcal{D} = \sum x_{ij}$
 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

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

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

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

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

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

Deriving $\hat{\mu}_{MLE}$

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

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

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

equating it to 0 gives

$$\Rightarrow \sum_{i=1}^N \sum^T (x^{(i)} - \bar{x}) = 0$$

$$\Rightarrow \sum_{i=1}^N (x^{(i)} - \bar{x}) = 0 \quad \left(\begin{array}{l} \text{multiply by } \Sigma \\ \text{on both sides} \end{array} \right)$$

$$\Delta \sum \Sigma^T = I \rightarrow \text{identity matrix.}$$

$$\Delta \sum \cdot 0 = 0 \rightarrow \text{null matrix.}$$

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

Ans.

Deriving Σ_{MLE} .

Deriving MLE for covariance matrix Σ , 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 \leftarrow it remains same.

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

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

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

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. Σ & using

$$|\Sigma^{-1}| = |\Sigma|^{-1} \text{ we get}$$

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

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

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

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

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

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

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

unknown parameter $\theta := \mu$

Q.5.

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

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

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

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

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.

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

$\because f_X(D)$ is independent of unknown

mean parameter μ

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

MAP means maximum a posteriori

mode $\rightarrow \mu$ at which $p(\mu | D)$

then it's maximum value:

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

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

Assuming samples $\mathcal{D} = \{x_i\}_{i=1}^N$ are ~~sample~~ i.i.d.

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

~~log likelihood~~ Prior

taking log on both sides gives.

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

lets call this $l(\theta)$

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

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

$$+ \log\left(\frac{1}{\sqrt{2\pi}\sigma_0}\right) - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2$$

taking differentiation on both side w.r.t. μ

~~green~~ and set it to zero given critical point

$$\frac{\partial l(\mu)}{\partial \mu} = 0 \Rightarrow \left[\sum_{i=1}^N \frac{1}{\sigma^2} (x_i - \mu) \right] - \frac{1}{\sigma_0^2} (\mu - \mu_0) = 0$$

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

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

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

So density function using parameters corresponding to mode of posterior distribution is

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

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

② Bayesian estimator corresponding to mean of posterior distribution

~~Exercise~~ ~~Refrigerator~~

Again by changing Using (7) from previous part

$$p(\mu | \mathcal{D}) = \frac{f_X(\mathcal{D} | \mu) \cdot b(\mu)}{f_X(\mathcal{D})} \quad f_X(\mathcal{D}) \sim \text{independent of } \mu$$

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

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

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

↳ say

Assuming samples are iid. we get ..

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

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

$$\propto e^{-\frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^2}} = e^{-\frac{(\mu - \mu_0)^2}{2\sigma^2}}$$

proportional symbol

\propto

$$-\sum_{i=1}^N (x_i - \mu)^2 \sigma_0^{-2} = -\sigma^2 (\mu - \mu_0)^2$$

$$\propto e^{-\frac{(\mu - \mu_0)^2}{2\sigma^2 \sigma_0^2}}$$

After simplification we get

$$\propto \exp \left\{ -\mu^2 \left(\sum_{i=1}^N \sigma_0^{-2} + \sigma^{-2} \right) + 2\mu \left[\left(\sum_{i=1}^N x_i \sigma_0^{-2} \right) + \mu_0 \sigma^2 \right] - \left(\left[\sum_{i=1}^N x_i^2 \sigma_0^{-2} \right] + \mu_0^2 \sigma^2 \right) \right\}$$

exp

$$\propto \exp \left\{ -\mu^2 + 2\mu \left[\frac{\sum x_i \sigma_0^{-2} + \mu_0 \sigma^2}{\sum \sigma_0^{-2} + \sigma^2} \right] - \left[\frac{\sum x_i^2 \sigma_0^{-2}}{\sum \sigma_0^{-2} + \sigma^2} \right] \right\}$$

Q. If $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$

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

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

then

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

$$\propto \exp \left(\frac{-(\mu - \mu_1)^2}{2\sigma_1^2} \right) \cdot \exp \left(\frac{C_1 + \mu_1^2}{2\sigma_1^2} \right)$$

also constant.

(by adding & subtracting μ_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 μ_1 & σ_1 , add a little

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

@ $1 + \frac{\sigma^2}{\sigma_0^2}$

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

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

and Variance σ_1^2 is

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

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

I guess we have 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 \approx

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

Key observation:

a) $N \rightarrow \infty$.

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

$\xrightarrow{\text{mean posterior estimation}}$ \downarrow 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 μ_{mean} is as defined above.

Q6. Devise EM steps for

car. Mixture of gaussians

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

Conditional distribution

$$(x_i | z_i = k) \sim N(\mu_k, \sigma_k^2)$$

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

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

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

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

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

Aim: find maximum likelihood estimates: π_j, μ_j & σ_j^2
 given a dataset $\sum_{i=1}^N x_i = x_i$

Optimization Problem

Likelihood function $L(\theta)$ is

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

Log likelihood will be

for notational convenience.

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

MLE $\hat{\theta}_{MLE} = \operatorname{argmax} l(\theta)$ s.t. $\sum_{j=1}^K \pi_j = 1$

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

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

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

$$\frac{\partial L(\theta, \lambda)}{\partial \sigma_j^2} = 0 \Rightarrow \left[\sum_{i=1}^N \frac{\pi_j f_i}{\sum_{j=1}^K \pi_j f_j} \cdot \left(-\frac{1}{2} f_j + f_j \cdot \left(\frac{-2}{\sigma_j^2} \right) \left(\frac{-1}{2} (x_i - \mu_j)^2 \right) \right) \right] = 0 \quad \text{--- (3)}$$

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.

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

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

$$③ \Rightarrow \sum_{i=1}^N \frac{\gamma_{ij}}{\pi_j} - \lambda = 0 \Rightarrow \sum_{i=1}^N \gamma_{ij} - \lambda \pi_j = 0$$

taking sum over j, on both sides gives.

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

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

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

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

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

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

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

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

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

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

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

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.

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

Aim: Find maximum likelihood estimator log π_j 's, p_j 's

given data $\sum x_i = n$ $\sum_{i=1}^N x_i$

likelihood function $L(\theta)$ is

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

log likelihood will be:

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

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

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

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

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

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

$$\frac{\partial L(\theta, \lambda)}{\partial p_j} = \frac{\sum \pi_j}{\pi_1 f_1 + \dots + \pi_K f_K} \cdot \left(x_i p_j^{x_i-1} (1-p_j)^{1-x_i} + (1-x_i) p_j^{x_i} (1-p_j)^{1-x_i} \right) = 0$$

$$\sum_{i=1}^N \frac{\pi_j b_j x_i (1-b_j)^{1-x_i}}{\pi_1 f_1 + \dots + \pi_k f_k} \left\{ x_i b_j^{-1} - (1-x_i)(1-b_j)^{-1} \right\} = 0$$

$$\sum_{i=1}^N \frac{\pi_j f_i}{\pi_1 f_1 + \dots + \pi_k f_k} \left\{ \cancel{x_i b_j^{-1} - (1-x_i)(1-b_j)^{-1}} \right\} = 0$$

$$\Rightarrow \sum_{i=1}^N \left[\gamma_{ij} = \frac{\pi_j f_i}{\pi_1 f_1 + \dots + \pi_k f_k} \right] \text{ from}$$

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

$$\hat{b}_j = \frac{\sum_{i=1}^N \gamma_{ij} x_i}{\sum_{i=1}^N \gamma_{ij}} \rightarrow \textcircled{Q}$$

From Q we also have

$$\sum_{i=1}^N \frac{f_i}{\pi_1 f_1 + \dots + \pi_k f_k} - \lambda = 0 \Rightarrow \sum_{i=1}^N \frac{\pi_j \gamma_{ij}}{\pi_j} - \lambda = 0$$

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

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

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

so $\hat{\pi}_j = \sum_{i=1}^N \gamma_{ij} \rightarrow \boxed{\hat{\pi}_j = \frac{\sum_{i=1}^N \gamma_{ij}}{N}}$

EM Algorithm

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

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

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

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

loop till convergence

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

$$f_{\mathbf{x}}(\mathbf{x} | b) = b^{x_i} (1-b)^{1-x_i}$$

given conjugate prior to be Beta distribution
i.e.

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

by bayes rule

$$f(b | D) = \frac{f_{\mathbf{x}}(D | b) \cdot g(b)}{f_{\mathbf{x}}(D)}$$

$f_{\mathbf{x}}(D)$ is independent of b . and is just a scaling constant.

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

$$b_{\text{MAP}} = \arg \max_b f_{\mathbf{x}}(\mathbf{x} | D)$$

$$= \arg \max_b \left(\prod_{i=1}^N f(x_i | b) \right) \cdot g(b)$$

taking log won't change argmax so,

$$b_{\text{MAP}} = \arg \max_b \sum_{i=1}^N \log(f(x_i | b)) + \log g(b)$$

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

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

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

By bayes rule.

Take differentiation w.r.t. p on both sides

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

"On simplification we get -"

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

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

Ans

Bayes estimate of

Q-8

Multinomial distribution

multinomial distribution.

~~D = {X_i}^N~~

$$\text{S.t. } \sum_{i=1}^K X_i = m$$

$$\sum_{i=1}^K x_{ij} = m$$

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

$$\text{where } Y_i = \{X_{ij}\}_{j=1}^{n_k}$$

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

then probability mass function of multinomial distribution is

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

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

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

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

Given Conjugate prior is dirichlet distribution

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

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

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

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

MAP estimate of multinomial distribution

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

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

Now we know that this is equivalent to

$$\vec{b}_{\text{MAP}} = \underset{b_1, b_2, \dots, b_K}{\operatorname{argmax}} \left(\sum_{i=1}^N \log (f(x_1^{(i)}, x_2^{(i)}, \dots, x_K^{(i)}) b_1, \dots, b_K) + \log g(b_1, \dots, b_K) \right)$$

Using Lagrange's multiplier

$$\Rightarrow \text{let } l(b_1, b_2, \dots, b_K) = \sum_{i=1}^N \log (f(x_1^{(i)}, x_2^{(i)}, \dots, x_K^{(i)}) b_1, \dots, b_K) + \log (g(b_1, \dots, b_K))$$

$$\begin{aligned} &= \sum_{i=1}^N \log (c_i) + \sum_{i=1}^N \sum_{j=1}^K x_j^{(i)} \log b_j + \log (\text{constant}) + \\ &\quad \sum_{j=1}^K b_j (\alpha_j - 1) \log b_j + \lambda (1 - \sum_{j=1}^K b_j) \end{aligned}$$

Now

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

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

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

- take summation on both sides w.r.t. j

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

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

↓
sum

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

Substitute (2) in (1) we get.

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

$$mN + \sum_{j=1}^k (d_j - 1)$$

for $j = 1, 2, \dots, k$.

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

A required MAP estimation of multinomial distribution.

Q.9

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 ^{unknown} parameters

$D = \{X_i\}$, $q \sim$ distribution of Z

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

$$= \underset{\Theta}{\operatorname{argmax}} \frac{p(\Theta, D)}{p(D)} \cdot \underbrace{p(D)}_{\text{prior distribution}}$$

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

($\because \log$ is increasing function).

$$= \underset{\Theta}{\operatorname{argmax}} \log p(\Theta, D) - \underset{\Theta}{\operatorname{argmax}} \log p(D)$$

independent of Θ

$$= \underset{\Theta}{\operatorname{argmax}} \log p(\Theta, D)$$

$$= \underset{\Theta}{\operatorname{argmax}} \log [p(D|\Theta) \cdot p(\Theta)]$$

likelihood

$$= \underset{\Theta}{\operatorname{argmax}} (\log p(D|\Theta) + \log p(\Theta))$$

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 or note z_i is hidden variable.)
 $\sum_c p(x_i; z_i=c | \theta) = p(x_i | \theta)$

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

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

s.t.

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

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

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

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

$\Rightarrow l(\theta) \geq l(\theta, q)$ decreasing (~~increasing~~) $\& q$ of $q(z_i=c) > 0$

↳ just a notation $\& \sum_c q(z_i=c) = 1$

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

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

for any q , since

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

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

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

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

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

↳ 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 E-step find q^*

$$\text{So } l(\theta) + l(\theta, q^*) = 0 \text{ when }$$

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

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

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

$$\Rightarrow q^*(z_i) := p(z_i | x_i, \theta) \rightarrow \text{independent of prior } b(\theta)$$

$$\text{So } \underline{\text{E-step}} \quad q^*(z_i) = p(z_i | x_i, \theta)$$

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

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

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

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

NOTE
 θ^* depends on prior $b(\theta)$

* proof of correctness -

$$\ell(\theta^{(old)}) \in \sum_{i=1}^N \sum_{c} q^*(z_i=c) \ln$$

$$\begin{aligned} \therefore \ell(\theta^{(old)}) &= \ell(\theta^{(old)}, q^*) + \sum_{i=1}^N \text{KL}(q^*(z_i) \parallel \\ &\quad \underbrace{P(z_i | x_i, \theta^{(old)})}_{\text{---}}) \\ &\quad (\because q^*(z_i) := p(z_i | x_i, \theta^{(old)})) \end{aligned}$$

$$\begin{aligned} \therefore \ell(\theta^{(old)}) &= \ell(\theta^{(old)}, q^*) \\ &= \sum_{i=1}^N \sum_c q^*(z_i=c) \log \left(\frac{P(x_i, z_i=c)}{q^*(z_i=c)} \right) + \log(b(\theta)) \end{aligned}$$

$$\begin{aligned} \ell(\theta^{(old)}) &= \sum_{i=1}^N \sum_c q^*(z_i=c) \log(P(x_i, z_i=c)) + \log(b(\theta)) \\ &\quad - \sum_{i=1}^N \sum_c \log(q^*(z_i=c)) \end{aligned}$$

$$\begin{aligned} \ell(\theta^{(new)}) &= \ell(\theta^{(new)}, q^*) + \sum_{i=1}^N \text{KL}(q^*(z_i) \parallel \\ &\quad \underbrace{P(z_i | x_i, \theta^{(new)})}_{\geq 0}) \\ &= \sum_{i=1}^N \sum_c q^*(z_i=c) \ln(P(x_i, z_i=c)) - \sum_{i=1}^N \sum_c q^*(z_i=c) \\ &\quad + \log(b(\theta)) + \text{KL}(q^*(z_i) \parallel P(z_i | x_i, \theta^{(new)})) \end{aligned}$$

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

$$\therefore \ell(\theta^{(new)}) \geq \ell(\theta^{(old)}) \Rightarrow \underline{H \cdot P}$$

END OF ASSIGNMENT

X X X

(cont'd) - 1)

(cont'd) - 2)

(cont'd) - 3)

(cont'd) - 4)

(cont'd) - 5)