

Advanced Machine Learning: Final Exam

Section A

A1.) Model 1 is too simple because it is a linear model. Compared to that the data points to be fitted are generated by a function of the 4<sup>th</sup> power. Thus, model 1 will end up underfitting on the data while model 2 will fit better to the data.

Thus, model 1 will have a higher bias than model 2

Ans: (B) Higher

A2.) For Bayes classifier, no of params:

$$3 \times (2^{10}) = 3072$$

For Naive Bayes classifier:  $3 \times (10+1) = 33$

Ans: (D) 3072, 33

A3.) Since it is a One vs Rest approach, the no of classifiers needed = no of classes. But when classes = 2, It is binary classification, so 1 will do.  
 $\therefore 5 \& 1$

Ans: (A) 5, 1      (B) 5, 1

A4) Since, we want to carry out the Expectation-Maximization algorithm, therefore we want to maximize the log likelihood.

$$\therefore \ln p(x|\theta) = L(q, \theta) + KL(q||p)$$

~~Therefore, we want to~~

$$= \sum_z q(z) \ln \left\{ \frac{p(x, z|\theta)}{q(z)} \right\}$$

$$- \sum_z q(z) \ln \left\{ \frac{p(z|x, \theta)}{q(z)} \right\}$$

④ In order to maximise  $\ln p(x|\theta)$ :

- we need to maximise  $L(q|\theta)$
- we need to minimise  $KL(q||p)$

Ans: ~~(c) Maximise, minimize~~

Ans: (c) Maximise, minimize

A5.) The probability distribution of the latent variables is hard to find and calculate. That is why an approximate value is calculated in the EM algorithm. The probability distribution  $q$  handles this by estimating the posterior.

Ans: (B)  $p(z|x, \theta)$

### Section: B

B1.) In leave-one-out cross validation 1 sample point is used as validation while remaining  $n-1$  samples are used for training. Thus, it is a special case of k-fold cross validation where  $k$  is no. of training samples.

Ans: True

B2.) Since, the test set data has been used to select model parameters, the test set error is biased over the training set error.

Ans: False

B3.) During gradient-descent for a logistic regression model, it may get stuck at a local minima as well.

Ans: False



B4.)

The convolution operations are linear transformations. If  $f_l$  is the non-linear activation function that introduce non-linearity. Without them, convolutional layers cannot do non-linear classification.

Ans: False

B5.)

When there is a small training set - cross validation is indeed used to train a model and tune the parameters.

Ans: True



## Section C

- 1)
- (i) Logistic regression is usually limited to a binary classification problem.  
But logistic regression can be extended to multi-class classification using a one-vs-rest approach. ~~where~~ In this approach, a separate classifier is trained for each class.
- (ii) A small training but a high validation error implies that the model is overfitting. Thus, it has a high variance and low bias. Thus, L2 regularization should be increased.
- (iii) Two ways to deal with non-linearly separable data using SVM are:
- Since SVM is a linear classifier, it introduces the concept of hyperplane which projects the otherwise non-linearly separable data into a higher dimension where it is linearly separable.
  - Combines soft margin and kernel trick to modify the decision boundary for non-linearly separable cases.

Q2)

6.5)  $f_x(x; \theta) = \theta x^{-\theta-2}$

$$L(\theta) = \prod_{i=1}^n f_x(x_i; \theta)$$
$$= \theta^n \prod_{i=1}^n x_i^{-\theta-2}$$

$$\ln(L(\theta)) = n \ln \theta + (-\theta - 2) \sum_{i=1}^n \ln x_i$$

$$\frac{\partial \ln(L)}{\partial \theta} = \frac{n}{\theta} + \sum_{i=1}^n \frac{1}{x_i}$$

$$\therefore \frac{\partial^2 \ln(L)}{\partial \theta^2} = -\frac{n}{\theta^2} < 0 \quad [\theta > 0]$$

$$\therefore \hat{\theta}_{MLE} = \frac{\partial \ln(L)}{\partial \theta} \quad \frac{\partial \ln(L)}{\partial \theta} = 0$$

$$\Rightarrow \hat{\theta}_{MLE} = \frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$$

$$Q2) p(\theta | x_1, x_2, \dots, x_n)$$

$$= \frac{p(x_1, x_2, \dots, x_n | \theta) \cdot p(\theta)}{p(x_1, x_2, \dots, x_n)} \rightarrow \text{normalization}$$

$$\therefore p(x_1, x_2, \dots, x_n | \theta) \cdot p(\theta)$$

$$= \theta^n \prod_{i=1}^n x_i^{-\theta-2} \lambda \theta^{-\lambda-2}$$

$$\Rightarrow \ln \{ p(x_1, x_2, \dots, x_n | \theta) \} = n \ln \theta + (-\theta - 2) \sum_{i=1}^n \ln x_i \\ \underbrace{+ \ln \lambda + (-\lambda - 2) \ln \theta}_{\ln p(\theta)}$$

$$\therefore \text{we want, } \frac{\partial p(\theta)}{\partial \theta} = 0$$

$$\Rightarrow \frac{n}{\theta} + \sum_{i=1}^n -\ln x_i - \frac{(-\lambda - 2)}{\theta} = 0$$

$$\frac{\partial^2 p(\theta)}{\partial \theta^2} = -\cancel{8} \quad \frac{n}{\theta^2} + \frac{(-\lambda - 2)}{\theta^2}$$

$$= -\frac{(n - \lambda - 2)}{\theta^2} < 0$$

$$\therefore \hat{\theta}_{MAP} = \frac{n-h-2}{\sum_{i=1}^n \ln x_i}$$

- Q3) a) The likelihood function  $L(\theta)$ , is defined which is the probability distribution of a set of observations ( $x$ ) ; given parameters ( $\theta$ )
- b) Log of the likelihood is taken as it makes computation easier
- c) The set of parameters need to be specified as several unknown parameters are calculated from the data
- d) An optimization of the MLE is performed. Here, we differentiate to get maxima
- e) Finally, the log likelihood function is optimized using the parameter values achieved from previous step.

Q3)

(a) Naive Bayes assumes all the features are independent, and have an equal contribution in predicting the class value. In this dataset there are 3 features, hence 3 independent parameters.

(b) Estimation of the parameter values

Size	Yes	No	$P(\text{Yes})$	$P(\text{No})$
small	1	3	$\frac{1}{4}$	$\frac{3}{6} = \frac{1}{2}$
large	3	3	$\frac{3}{4}$	$\frac{3}{6} = \frac{1}{2}$
Total	4	6	1	1

Shape	Yes	No	$P(\text{Yes})$	$P(\text{No})$
Circle	3	2	$\frac{3}{4}$	$\frac{2}{6} = \frac{1}{3}$
Irregular	1	4	$\frac{1}{4}$	$\frac{4}{6} = \frac{2}{3}$
Total	4	6	1	1

Color	Yes	No	$P(\text{Yes})$	$P(\text{No})$
Red	4	2	$\frac{1}{2}$	$\frac{1}{6}$
Green	0	5	0	$\frac{5}{6}$
Total	4	6	1	1

Good Apple

	Yes	4
	No	6
Total		10

$$P(\text{Yes}) = \frac{4}{10} = \frac{2}{5}$$

$$P(\text{No}) = \frac{6}{10} = \frac{3}{5}$$

(Q3) New feature  $x = \{\text{Small, Red, Round}\}$

$$P(y=\text{No}|x) = \frac{P(\text{small}|\text{No}) \cdot P(\text{Red}|\text{No}) \cdot P(\text{Round}|\text{No}) \cdot P(\text{No})}{P(x)}$$

$$= \frac{3}{6} \times \frac{1}{6} \times \frac{3}{6} \times \frac{6}{10}$$

$$= 0.016$$

$$P(y=\text{Yes}|x) = \frac{P(\text{small}|\text{Yes}) \cdot P(\text{red}|\text{Yes}) \cdot P(\text{round}|\text{Yes}) \cdot P(\text{Yes})}{P(x)}$$

$$= \frac{1}{4} \times \frac{4}{4} \times \frac{3}{4} \times \frac{4}{10}$$

$$= 0.075$$

In order to convert it to a probability value

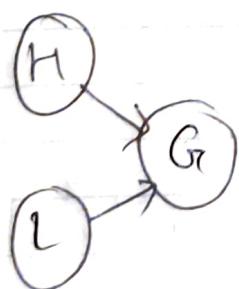
$$p(\text{Yes} | x) = \frac{0.075}{0.075 + 0.016}$$
$$= 0.824$$

$$p(\text{No} | x) = \frac{0.016}{0.075 + 0.016}$$
$$= 0.176$$

$\because p(\text{Yes} | x) > p(\text{No} | x)$ , the Naive Bayes classifies  $y = \text{Yes}$  for  $x$

4)  $G_1 = \text{Going out} ; H = \text{Homework} ; L = \text{Lawnmower}$

(a)



$$p(H) p(L) p(G_1 | H, L)$$
$$= p(G_1, H, L)$$

$\because$  Going out depends on having homework and if Lawmower is available

(Q2) Let the following hold:

$$P(G, H, L) = P_1$$

$$P(\bar{G}, H, L) = P_6$$

$$P(\bar{G}, \bar{H}, L) = P_2$$

$$P(\bar{G}, \bar{H}, \bar{L}) = P_7$$

$$P(G, \bar{H}, L) = P_3$$

$$P(G, H, \bar{L}) = P_4$$

$$P(\bar{G}, \bar{H}, \bar{L}) = P_5$$

$$\therefore P(\bar{G}, \bar{H}, \bar{L}) = 1 - (P_1 + P_2 + P_3 + P_4 + P_5 + P_6 + P_7)$$

Let:  $P(H) = P_h$        $P(\bar{H}) = 1 - P_h$

$P(L) = P_e$        $P(\bar{L}) = 1 - P_e$

$$\therefore P(G, H, L) = P(H) P(L) P(G|H, L)$$

$$P_1 = P_h \times P_e \times P_g$$

then,  $P(\bar{G}, H, L) = P(\bar{H}) \times P(L) P(\bar{G}|H, L)$

$$P_2 = P_h \times P_e \times (1 - P_g)$$

$$\Rightarrow P_2 = P_h \times P_e - P_1$$

$P_1$  &  $P_2$  are interrelated

$$P(G, \bar{H}, L) = P(\bar{H}) P(L) P(G | \bar{H}, L)$$

$$P_8 = (1 - P_h) P_L P_{G''}$$

$$P(\bar{G}, \bar{H}, L) = P(\bar{H}) P(L) P(\bar{G} | \bar{H}, L)$$

$$P_5 = (1 - P_h) P_L (1 - P_{G''})$$

$$P_5 = (1 - P_h) P_L - P_3$$

$P_5$  &  $P_3$  are also related

Similarly,  $P_4$  &  $P_6$  are also related.

$$\therefore P(\bar{G}, \bar{H}, \bar{L}) = 1 - (P_1 + P_2 + P_3 + P_4 + P_5 + P_6 + P_7)$$

$\therefore (P_1 \& P_2), (P_3 \& P_5), \& (P_4 \& P_6)$   
are related

Thus, it depends on 6 parameters:

$p_1/p_2$ ;  $p_2/p_4$ ;  $p_3/p_5$ ;  $p_7$ ,  $p_h$  and  $P_L$

- (Q3) After G has been observed, H and L do not depend on each other and are thus independent-

Q4)

$$\begin{array}{c}
 \text{A} \rightarrow \text{L} \rightarrow \text{G} \\
 \text{H} \rightarrow \text{G}
 \end{array}
 \quad p(G, H, L, A) = p(H)p(A)p(L|A)p(G|H, L)$$

Thus, if  $L$  is given then  $A$  &  $G$  are independent of each other

Q5)

$$\text{A} \rightarrow \text{L} \rightarrow \text{G}$$

$$p(G, L, A) = p(A)p(L|A)p(G|L)$$

$$p(G, L, A) = q_1 \quad p(\bar{G}, L, A) = q_2 \quad p(G, \bar{L}, A) = q_3$$

~~$p(\bar{G}, \bar{L}, A)$~~

$$p(G, L, \bar{A}) = q_4 \quad p(\bar{G}, \bar{L}, A) = q_5 \quad p(\bar{G}, L, \bar{A}) = q_6$$

$$p(G, \bar{L}, \bar{A}) = q_7$$

$$\begin{aligned}
 \text{let: } p(A) &= p_a & p(L|\bar{A}) &= p_e \\
 && p(L|A) &= p_g
 \end{aligned}$$

$$p(G, L, A) = p(A)p(L|A)p(G|L)$$

$$q_1 = p_a p_g p_g \quad (\text{assume})$$

$$p(\bar{G}, L, A) = p(A)p(L|A)p(\bar{G}|L)$$

$$q_2 = p_a p_e (1 - p_g)$$

$$\Rightarrow q_2 = PaPe - q_1 \quad - (i)$$

$\therefore q_1$  and  $q_2$  are related

$$P(G, \bar{L}, A) = P(A) P(\bar{L}|A) P(G|\bar{L})$$

$$q_3 = P_a (\bar{J} - P_e) P_{g'} \text{ (assume)}$$

$$P(\bar{G}, \bar{L}, A) = P(A) P(\bar{L}|A) P(\bar{G}|\bar{L})$$

$$q_5 = P_a (\bar{J} - P_e) (\bar{J} - P_{g''})$$

$$\Rightarrow q_5 = P_a (\bar{J} - P_e) - q_3$$

$\therefore q_3$  and  $q_5$  are also related

Also,

$$P(A, L, \bar{A}) = P(\bar{A}) P(L|\bar{A}) P(G|L)$$

$$q_4 = (\bar{J} - P_a) P_e' P_{g''} \text{ (assume)}$$

$$P(\bar{G}, L, \bar{A}) = P(\bar{A}) P(L|\bar{A}) P(\bar{G}|\bar{L})$$

$$q_6 = (\bar{J} - P_a) P_e' (\bar{J} - P_{g''})$$

$$\Rightarrow q_6 = (\bar{J} - P_a) P_e' - q_4$$

$\therefore q_4$  and  $q_6$  are also related.

Thus for the joint probability  $P(\bar{G}, \bar{L}, \bar{A})$   
we will need  $q_3/q_2$ ;  $q_5/q_5$ ;  $q_4/q_6$

; q7; Pa; Pe and Pe'

5)

(Q) Size of image :  $3 \times 32 \times 32$

This is in order of  $(C \times H \times W)$

Originally image is :  $32 \times 32 \times 3$   
 $(H \times W \times C)$

Conv2D : input channel  $\rightarrow 3$   
output channel  $\rightarrow 6$   
Kernel  $\rightarrow 3 \times 3$       Pooling  
Kernel  $\rightarrow 2 \times 2$   
Stride  $\rightarrow 2$

Size of image after convolution:

$$\text{if } H' = W' = \text{ref } \frac{H - K + 2P}{S} + 1 \\ = 32 - 3 + 1 \\ = 30$$

$\therefore$  output channel = 6

$\therefore$  size of image =  $(30 \times 30 \times 6)$

Size after max pooling :

$$\text{size} = \frac{H - K + 1}{S} = \frac{30 - 2 + 1}{2} \\ = 15$$

∴ Final image size =  $(15 \times 15 \times 6)$

Q2) Max pooling outputs max value in each window position, while average pooling will give an average of the features in the window.

Size of feature map:  $4 \times 4$

Max pooling:

∴ Output shape  $\rightarrow 2 \times 2$

Kernel  $\rightarrow 2 \times 2$

Stride  $\rightarrow 2$

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

Max Pooling:

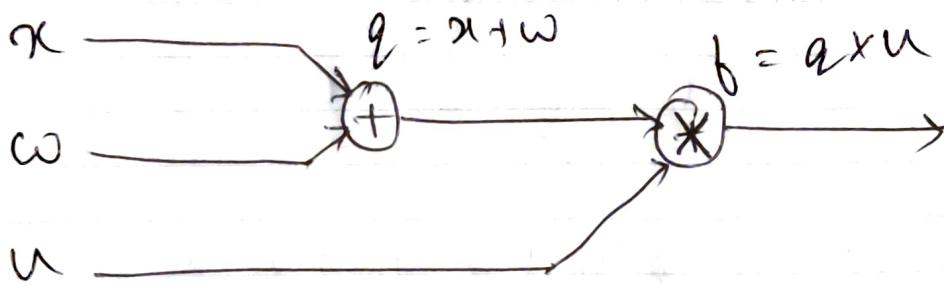
6	8
3	4

Averg. Pooling

3	5
2	2

B3)  $x = -2 \quad \omega = 5 \quad u = -4$  (initial values)

~~Q = x + \omega~~



$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \times \frac{\partial q}{\partial x} = \cancel{\frac{\partial (qxu)}{\partial q}} \times \frac{\partial (x+\omega)}{\partial x}$$

$$= u \times 1 = -4$$

$$\frac{\partial f}{\partial \omega} = \frac{\partial f}{\partial q} \times \frac{\partial q}{\partial \omega} = \cancel{\frac{\partial (qxu)}{\partial q}} \times \frac{\partial (x+\omega)}{\partial \omega}$$

$$= u \times 1 = -4$$

$$\frac{\partial f}{\partial u} = \cancel{\frac{\partial (qxu)}{\partial u}} = q = (x+\omega) = 3$$

6)

Q1)

We need to choose 3 clusters from the 2D space where there are 6 points

$${}^6C_3 = \frac{6 \times 5 \times 4}{3 \times 2} = 20$$

	<u>3 partition</u>	stable?	An ex of 3 start config
i)	{a,b}, {c,f}, {d,e}	yes	a,c,d
ii)	{a,b,e}, {c,d}, {f}	no	none
iii)	{a,d}, {b,c}, {e,f}	yes	a,b,e
iv)	{a,b,d}, {e}, {c,f}	no	none
v)	{a,d,e}, {b,c}, {f}	yes	d,c,f

Explaining using (iii) as an example:

Let the initial cluster centroids be a, b, e

- 1) For a(0,0)  $\rightarrow$  nearest centroid a(0,0)  
 For b(8,0)  $\rightarrow$  nearest centroid b(8,0)  
 For c(16,0)  $\rightarrow$  nearest centroid b(8,0)  
 For d(0,6)  $\rightarrow$  nearest centroid a(0,0)  
 For e(8,0)  $\rightarrow$  nearest centroid e(8,0)  
 For f(16,6)  $\rightarrow$  nearest centroid e(8,6)

2) For centroid movement:

a = centroid for a and d  $\rightarrow$  cluster 1

$$(x_{\text{cent}})_1 = \frac{(a_{\text{cent}})_x + (d_{\text{cent}})_x}{2} = 0$$

$$(y_{\text{cent}})_1 = \frac{(a_{\text{cent}})_y + (d_{\text{cent}})_y}{2} = 3$$

$\therefore$  centroid is at  $(0, 3)$ , hence converged

b = centroid for b & c = cluster 2

$$(x_{\text{cent}})_2 = \frac{(b_{\text{cent}})_x + (c_{\text{cent}})_x}{2} = 12$$

$$(y_{\text{cent}})_2 = \frac{(b_{\text{cent}})_y + (c_{\text{cent}})_y}{2} = 0$$

$\therefore$  centroid is at  $(12, 0)$  and so it converged

e = centroid for e & f  $\rightarrow$  cluster 3

$$(x_{\text{cent}})_3 = \frac{(e_{\text{cent}})_x + (f_{\text{cent}})_x}{2}$$

$$= 12$$

$$(y_{\text{cent}})_3 = \frac{(e_{\text{cent}})_y + (f_{\text{cent}})_y}{2}$$

$$= 6$$

$\therefore$  centroid is at  $(12, 6)$  & converged

New centroids:

$$(0, 3); (12, 0); (12, 6)$$

For	a $(0, 0) \rightarrow$ nearest centroid	$(0, 3)$
b	$(8, 0) \rightarrow$ nearest centroid	$(12, 0)$
c	$(16, 0) \rightarrow$ nearest centroid	$(12, 0)$
d	$(0, 6) \rightarrow$ nearest centroid	$(0, 3)$
e	$(8, 6) \rightarrow$ nearest centroid	$(12, 6)$
f	$(16, 6) \rightarrow$ nearest centroid	$(12, 6)$

$\because$  clusters remain the same, re-calculating will lead to same clusters.

Thus, it has converged & is stable

7)

(Q1)

given,

$$R: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^n$$

$R$  is valid if there is a feature vector  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^m$  s.t-

$$R(x, x') = \phi(x)^T \phi(x')$$

$K_1(x, x')$ ,  $K_2(x, x')$   $\rightarrow$  valid kernels

For  $R(x, x')$  to be a valid kernel the gram matrix  $K_m$  must be +ve semidefinite for  $\{x_n\}$  where elements in  $K_m$  are  $R(x_n, x_m)$

For  $K_m$  to be positive semi-definite

$$y^T K_m y \geq 0 \quad (\rightarrow y \text{ is any vector})$$

$K_{m_1} \rightarrow$  gram matrix of  $R_1(x, x')$   
 $K_{m_2} \rightarrow$  gram matrix of  $R_2(x, x')$

$\therefore K_{m_1}$  &  $K_{m_2}$  are +ve semi-definite

$$y^T K_{m_1} y \geq 0$$

$$y^T K_{m_2} y \geq 0$$

$\therefore$  Gram matrix of  $k_1(x, x')$  &  $k_2(x, x')$

$$= Km_1 + Km_2$$

$$= y^T (Km_1 + Km_2) y$$

$$= y^T Km_1 y + y^T Km_2 y \geq 0$$

$\therefore K(x, x') = k_1(x, x') + k_2(x, x')$  is

a valid Kernel

(Q2)  $K_1(x, x') = f(x)^T f(x')$  } (assuming)  
 $K_2(x, x') = g(x)^T g(x')$

$$\begin{aligned} k_1(x, x') \cdot K_2(x, x') &= f(x)^T f(x') g(x)^T g(x') \\ &= \sum_{p=1}^X f_p(x) f_p(x') \sum_{q=1}^Y g_q(x) g_q(x') \\ &= \sum_{p=1}^X \sum_{q=1}^Y f_p(x) f_p(x') g_q(x) g_q(x') \\ &= \sum_{m=1}^M \phi(x) \phi(x') = \phi(x)^T \phi(x') \end{aligned}$$

Where;  $M = XY$

$\therefore K_1(x, x')$  &  $K_2(x, x')$  are valid kernels

(Q3)  $K_3(x, x') \rightarrow$  a valid kernel

also,  $K_3(x, x') = f(x)^T f(x')$

$f(x) \rightarrow$  feature vector

$B_x \Rightarrow$  kernel const.

$$B_0 K_3(x, x') = h(x)^T h(x') \quad (h(x) = \sqrt{B_0} f(x))$$

$B_0 K_3(x, x')$  is the scalar product of feature vectors. So, it is a valid kernel

Also,  $p(x) K_3(x, x') p(x')$

$$= p(x) f(x)^T f(x') p(x')$$

$$= \phi(x)^T \phi(x') \quad [\phi(x) = p(x) f(x)]$$

$\rightarrow$  a scalar product of coefficients of all orders in  $p(x)$  &  $f(x)$

$\therefore R(x, x') = p(x) K_3(x, x') p(x')$  is a valid kernel

$$Q4) e^x = 1 + x + \frac{x^2}{2!} + \dots + \frac{x^n}{n!} \quad (\text{Taylor series})$$

$$k_1(x, x') = f(x)^T f(x')$$

$$e^{R_2(x, x')} = 1 + \underbrace{f(x)^T f(x')}_{2!} + \frac{\underbrace{[f(x)^T f(x')]}_{n}}{n!}$$

$$+ \dots + \frac{\underbrace{[f(x)^T f(x')]}_{n!}}{n!}$$

$$[f(x)^T f(x')]^n = \underbrace{[f(x)^T f(x')]}_{n\text{-times}} x \dots$$

also,  $f(x)^T f(x')$  are valid kernels

$\therefore e^{R_2(x, x')}$  is a valid kernel

$$Q5) R_1(x, x') = \exp \left\{ -\frac{\|x - x'\|^2}{2\sigma^2} \right\}$$

$$\Rightarrow \|x - x'\|^2 = x^T x - 2x^T x' + (x')^T x'$$

we already know,  $e^{R_1(x, x')} \rightarrow$  valid kernel

$\therefore R_1(x, x') \rightarrow$  valid kernel

$$= \frac{1}{2\sigma^2} [x^T x - 2x^T x' + (x')^T x']$$

$$= e^{-\frac{x^T x}{2\sigma^2}} \cdot e^{-\frac{x^T z'}{\sigma^2}} \cdot e^{-\frac{(z')^T (x')}{2\sigma^2}}$$

$$= g(z) e^{k_z(z, z')} \cdot g(z')$$

where  $R_z(z, z') = B_0 x^T x$

~~$R_z(x, x') = \exp \left\{ -\frac{\|x - x'\|^2}{2\sigma^2} \right\}$~~  is  
a valid Kernel

8)  $p(x) = \sum_{k=1}^K \pi_k p(x|k)$

mean  $\{k_1, \dots, k_n\} = \{m_1, \dots, m_n\}$   
covariance  $\{ \sum_{j=1}^n \dots \sum_{n=1}^n \}$

(a)  $E(x) = \int_x x p(x) dx$

$$= \int_x x \sum_{k=1}^K \pi_k p(x|k) dx$$

$$= \int_x x \{ \pi_1 p(x|k=1) + \dots + \pi_K p(x|k=K) \} dx$$

$$= \pi_1 m_1 + \pi_2 m_2 + \dots + \pi_K m_K$$

$E(x) = \sum_{k=1}^K \pi_k m_k$

$$\text{Q2)} \quad \begin{aligned} E[X|K=k] &= \mu_k \\ \text{cov}[X|K=k] &= \sigma_k^2 \end{aligned}$$

$$\therefore \text{cov}[X_i, Y_j | K=k] = E(X_i Y_j | K=k) - E(X_i | K=k) E(Y_j | K=k)$$

$$\text{And, } E(X) = \sum_{k=1}^K \pi_k \mu_k$$

$$\therefore E(X_i Y_j) = \sum_{k=1}^K E(X_i Y_j | K=k) \cdot \pi_k$$

$$= \sum_{k=1}^K \pi_k [(\Sigma_k)_{ij} + E(X_i | K=k) + E(Y_j | K=k)]$$

$$= \sum_{k=1}^K \pi_k [(\Sigma_k)_{ij} + (\mu_k)_i + (\mu_k)_j]$$

$$= \underbrace{\sum_{k=1}^K (\mu_k)_i \pi_k}_{E(X)} + \underbrace{\sum_{k=1}^K (\mu'_k)_j \pi_k}_{E'(X)}$$

$$\therefore \text{cov}(X_i, Y_j) = E[\text{cov}(X_i, Y_j | z)]$$

$$+ \text{cov}(E(X_i | z), E(Y_j | z))$$

Thus, from law of co-variance

$$\text{Cov}[\kappa] = \sum_{k=1}^K \pi_k (\Sigma_k + M_k U_k^T) - E(\kappa) E^T(\kappa)$$

q)

(a)  $L(q) = \int q(z) \frac{\ln p(x, z)}{q(z)} dz$

$$KL(q \parallel p) = - \int q(z) \frac{\ln p(z|x)}{q(z)} dz$$

Product Rule:

$$L(q) = \int q(z) \ln \left\{ \frac{p(x|z) \cdot p(z)}{q(z)} \right\} dz$$

$$p(x|z) = \frac{p(z|x) \cdot p(x)}{p(z)}$$

$$\therefore L(q) = \int q(z) \ln \left\{ \frac{p(z|x) \cdot p(x)}{p(z) q(z)} \right\} dz$$

$$= \int q(z) \left\{ \ln \left( \frac{p(z|x)}{q(z)} \right) + \ln p(x) \right\} dz$$

$$= \int q(z) \frac{\ln p(z|x)}{q(z)} dz + \int \ln(p(x)) dz$$

$$L(q) = -k \lambda (q \| p) + \ln P(x)$$

$$\therefore \ln P(x) = k \lambda (q \| p) + L(q)$$

Hence proved

$$(82) \quad q(z) = \prod_{i=1}^M q_i(z_i)$$

$$KL(p \| q) = - \int p(z) \left[ \sum_{i=1}^M \ln q_i(z_i) \right] dz + c$$

$\therefore$  Optimising wrt  $q_K(z_K)$

$$KL(p \| q) = - \int p(z) \ln q_K(z_K) dz + \int p(z) \sum_{i \neq K} \ln q_i(z_i) dz + c$$

$$= - \int p(z) \ln q_K(z_K) dz + c$$

$$= - \int \ln q_K(z_K) \left[ \int p(z) \pi_{i \neq K} dz_K \right] dz$$

$$= - \int B_K(z_K) \ln q_K(z_K) dz_K + c$$

$$\text{Now, } B_K(z_K) = \int p(z) \pi_{i \neq K} dz_K$$

Now using Lagrange's multiplier: so the  $q_k(z_k)$  is:

$$\begin{aligned} q_k(z_k) &\in \mathcal{L}_f \\ -\int B_k(z_k) \ln q_k(z_k) dz_k + \lambda \int q_k(z_k) dz_k^{-1} \end{aligned}$$

Taking the derivative and equating to 0

$$\frac{-B_k(z_k)}{q_k(z_k)} + \lambda = 0$$

$$\Rightarrow \lambda q_k(z_k) = B_k(z_k)$$

Integrating  $\rightarrow$

$$\lambda = \int B_k(z_k) dz_k$$

$$= \int [ \sum_{i \neq k} p(z) \pi_i dz_i ] dz_k = 1$$

$$\therefore q_k(z_k) = \frac{\int p(z) \pi_i dz}{\int_{i \neq k} p(z) \pi_i dz}$$

10)

(a) The likelihood of a dataset -  $P(X|\theta)$  can be computed using the EM algorithm.

In the E step the current parameters are used to compute posterior distribution of latent variables  $P(z|x, \theta)$

In the M step the current parameters are used to compute the posterior distribution of latent variables

In the M step - the parameters are optimized

$$\theta^{\text{new}} = \arg\max Q(\theta, \theta^{\text{old}})$$

$$Q(\theta, \theta^{\text{old}}) = \sum_z P(z|x, \theta^{\text{old}}) \ln P(x, z|\theta)$$

$$\alpha(z_t) = P(x_1, \dots, x_t, z_t)$$

$$\beta(z_t) = P(x_{t+1}, \dots, x_T | z_t)$$

$$\xi(z_{t+1}, z_t) = P(z_{t+1}, z_t | x)$$

$$= \frac{\alpha(z_{t+1}) P(x_t | z_t) P(z_t | z_{t-1}) \beta(z_t)}{\sum_{z_t} \alpha(z_t)}$$

Also,

$$\pi_k = \frac{\cancel{\gamma(z_k)}}{\sum_{j=1}^K \gamma(z_j)}$$

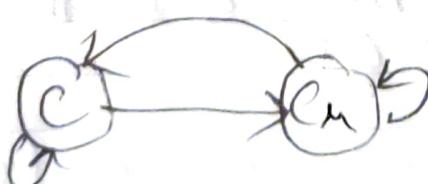
$$A_{jk} = \frac{\sum_{t=1}^T \mathbb{E}(z_{t-1}, j z_{tk})}{\sum_{m=1}^K \sum_{t=2}^T \mathbb{E}(z_{t-1}, j z_{tk})}$$

∴ From  $\mathbb{E}(z_{t-1}, j z_{tk})$ , if  $A_{jk} = 0$   
then  $\mathbb{E}(z_{t-1}, j z_{tk}) = 0$

Thus, no more updates in EM algorithm

Q2)

Chocolate lever = C  
Caramel lever = Cr



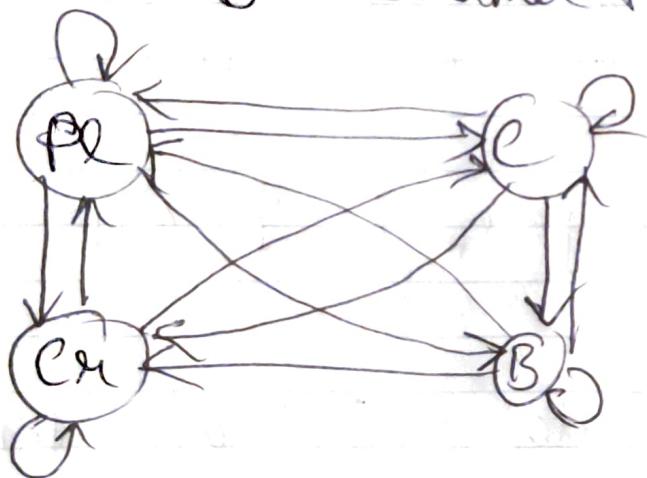
Prior probability as?

$$P(C) = P(Cr) = 0.5$$

Transition probability table :

$T \rightarrow$	C	Cr
e	0.7	0.3
Cr	0.3	0.7

Let :  $P_e =$  Plain  
 $B =$  Caramel + Chocolate (Both)



$\therefore$  levers are ind.

$$\therefore P(P) = P(C) = P(Cr) = P(B) = 0.25$$

State transition probability of both being flipped  $= (0.3)^2 = 0.09$

State transition probability of one being flipped

$$= 0.3 \times 0.7 = 0.21$$

State transition probability of neither being flipped

$$= 1 - 0.09 - 0.21 - 0.21 \\ = 0.49$$

$\therefore T =$

	PL	C	CM	B
PL	0.49	0.21	0.21	0.09
C	0.21	0.49	0.09	0.21
CM	0.21	0.09	0.49	0.21
B	0.09	0.21	0.21	0.49

Q3) Probability of the order:

$P\{ \text{Plain, Chocolate, Chocolate, Chocolate + Caramel} \}$

$$= P(C|PL, e, C, B)$$

$$= P(C|P) \cdot P(C|PL) \cdot P(C|e) \cdot P(B|C)$$

$$= 0.25 \times 0.21 \times 0.49 \times 0.21$$

$$= 0.0054$$