

Section A. Multiple Choice

A1) (B) Higher.

The set of data points follow a 4th order process. Therefore, we may observe that the Model 1 (1st order polynomial) underfits on the data while the Model 2 (9th order polynomial) overfits on the same data. This happens because the Model 1 is too simple to grasp the complexity present in the 4th order data, while the Model 2 is too complex and may pick up the noise variable which is present in the data along with the actual variables, overfitting to the training data. Regarding bias-variance, the Model 2 will have higher variance (too sensitive to noise present in the training data), while the Model 1 will have higher bias (unable to identify the key patterns in the training data because it is too simple).

A2) (C) 3071, 32.

All features are dependent on each other for a Bayes Classifier because we do not have the conditional independence assumption leading to a large number of parameters. The likelihood probabilities and prior probabilities add up to the total number of independent parameters.

In case of 3 classes, if we have prior probabilities of 2 classes, the prior for the 3rd class can be calculated (as their sum is 1). For 10 binary features and 3 classes, we need $3 * (2^{10} - 1) = 3069$ features (subtracting 1 because we can find the last probability using the previous all probabilities).

In total we have 3071 independent features.

For a Naïve Bayes classifier, we assume features are independent. Therefore, for 3 classes and 10 binary features we have likelihood features $3 * 10 = 30$. Prior probabilities are the same, i.e. 2. In total 32 independent features are required for a Naïve Bayes Classifier.

A3) (B) 5, 1

With One-vs-rest approach, we will need 5 classifiers for 5 classes, each corresponding to one of the classes. The predicted class is given by the classifier which has the highest output value. For 2 classes, we will need only 1 classifier since it is binary classification task.

A4) (C) Maximize, minimize

The EM algorithm aims to maximize the likelihood of the data by maximizing the variational bound. To do this, the algorithm must minimize the KL divergence between the true posterior distribution and the approximating distribution. Ideally, if the approximating distribution represents the true posterior exactly, the KL divergence will be 0.

A5) (B) $p(Z|X, \theta)$

We use $q(Z)$ distribution to simplify the calculation of the posterior $p(Z|X, \theta)$ given X . We use the EM algorithm to try and minimize the distance between these two distributions while maximizing the log-likelihood.

Section B. True/False

B1) True.

K-fold cross validation is a method for evaluating a model by dividing the dataset into k subsets (folds) and training the model on k-1 subsets and evaluating it on the remaining subset. This is repeated for each fold, and the overall performance is the average of the performance on each fold. When k equals the number of samples in the dataset, this is called leave-one-out cross validation, because each sample is used as the validation set once and the rest of the samples are used as the training set.

B2) False.

Using the test set to select model parameters can bias the model's error. To get a true, unbiased estimate of the model's performance, we need to use three data splits: a training set to train the model, a validation set to choose the hyperparameters, and a test set to evaluate the model's final performance.

B3) False.

We try to find the minimum of a cost function using gradient descent iteratively. The convergence of this algorithm depends on the function of which we are trying to find the global minimum, since there may be a case where the algorithm may get stuck in a local minimum.

B4) False.

We need activation functions to inject non-linearity in the classifier and we typically use non-linear transformations. Convolutional layers stacked onto each other will only transform the data linearly as they are just sets of filters that apply convolution operation on input and these filters are learned through training. Therefore, it is not possible to construct a non-linear classification model using just convolutional layers.

B5) True.

The model will have less information to train on if the training dataset size is small. If the validation dataset is also small, the error produced on it is less reliable. If the dataset is small, both training and validation sets are small. To avoid this, we use cross-validation, where we create folds from the dataset, using one fold for validation and rest for training, and repeating this for as many folds we have. We average the model performances and get the final performance. In this approach, the model gets to train and validate on a larger number of samples. This approach is computation intensive, and therefore not suggested for a large dataset.

B6) True.

For a Naïve Bayes classifier:

$$c = \operatorname{argmax}_c \left(p(C = C_j) \prod_i p(X = X_i | C = C_j) \right)$$

In a first-order Markov chain: $p(C = C_j)$ is $p(C = C_j | C_{j-1})$

But if we remove the conditional dependence between sequential observations, it becomes $p(C = C_j)$. Therefore, both first-order Markov chain and the Naïve Bayes classifier become the same.

B7) False.

K-means does not follow the same approach to determine cluster locations like GMM, and therefore is not a special application of GMM. In k-means, data points are assigned to the nearest centroid using some distance metric, but in GMM data points are assigned with likelihood of being generated by the current set of clusters.

Section C. Written Questions

1: Logistic Regression/Support Vector Machines

Q1) Logistic regression is inherently designed for binary classification where we model the target using a binomial probability distribution. In its original state, logistic regression cannot be used for multi-class classification and requires tweaking the methodology to support multi-class classification. There are two approaches to apply logistic regression for multi-class classification problems:

1. A popular approach is splitting the multi-class classification problem into multiple binary classification tasks using techniques like **one-vs-rest** and one-vs-one wrapper models.
In the one-vs-rest technique, taking an example of 5-class classification problem we will train 5 different logistic regression classifiers.
In the training phase, we will treat input data for a specific class as positive samples ($y=1$), and input data for rest all the classes as negative samples ($y=0$). We will repeat this for all the classes.
In the prediction phase, we will compute the probabilities for all the classes using all the classifiers and the final prediction function will choose the class for which the corresponding classifier gives the highest probability for the class label on the input example.
2. The second approach is using **multinomial logistic regression** based on multinomial probability in order to predict the class probability directly. This is done by changing the loss function to train the model from log-loss to cross-entropy loss, and changing the output from a single probability to an array of probability values (one for each class).

Softmax Regression (Multinomial Logistic Regression)

In Softmax Regression, the probability that a data point belongs to each class is calculated by:

$$\begin{bmatrix} P(y=1|x;\theta) \\ P(y=2|x;\theta) \\ \vdots \\ P(y=K|x;\theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^K \exp(\theta^{(j)T}x)} \begin{bmatrix} \exp(\theta^{(1)T}x) \\ \exp(\theta^{(2)T}x) \\ \vdots \\ \exp(\theta^{(K)T}x) \end{bmatrix}$$

Normalizes the probabilities so they sum to 1.

→ Predict the class with highest probability

Q2) If a logistic regression model achieves small training errors but large validation errors, it is likely overfitting the training data. In this case, increasing the L2 regularization would be appropriate because it would help to reduce the model's complexity and prevent overfitting. This is because L2 regularization adds a penalty to the model's weights, which encourages the model to use only the most important features for making predictions and discourages the use of unnecessary and potentially overfitting features. In terms of the bias-variance trade-off, increasing L2 regularization would help to reduce the model's variance (i.e., the amount by which the model's predictions vary depending on the training data used) and potentially increase its bias (i.e., the error due to the model's assumptions or limitations). This trade-off can be helpful in preventing overfitting and improving the model's generalization performance on unseen data.

Q3) We can deal with data that is not linearly separable using SVMs using two different methods:

- I) **Soft Margin:** In this approach, we still try to find a hyperplane that separates the data into two classes, but we tolerate misclassifications by introducing a margin. There would be two types of misclassifications after introducing a margin: first where a data point is on the incorrect side of the decision boundary but on the correct side of the margin (first diagram), and second where the data point is on the wrong side of the decision boundary as well as on the wrong side of the margin.

We try to adjust the search for a line for maximum margin and minimum misclassification on training set. This trade-off is an important hyperparameter for SVM and is controlled by the term 'C' in Sklearn. For large values of C, we get a hyperplane with a small margin that gets correctly classifies all training points. Conversely, when C has a very small value, the optimization problem finds a hyperplane with a large margin which may misclassify more training points. Therefore, training an SVM requires one to balance the importance of the regularization term with respect to the loss term using C.

- II) **Kernel Trick:** In this approach, we transform the linearly inseparable data by projecting it into a higher dimension and making it linearly separable there, using kernel functions and thus effectively find a non-linear decision boundary. Some of the most popular kernels are the Polynomial and Radial Basis Function(RBF).

2:Maximum Likelihood/Maximum A Posteriori Estimation

Q1) Given that the observations are i.i.d., we can estimate the value of the parameter θ using the maximum likelihood estimate (MLE). The likelihood function is the probability of observing the given data given a specific value of the parameter θ . In this case, the likelihood function is given by

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

The MLE is the value of θ that maximizes the likelihood function. Taking the natural log of the likelihood function, we can simplify the problem to finding the value of θ that maximizes the log-likelihood function:

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

Taking the derivative of the log-likelihood function with respect to θ and setting it equal to zero, we can find the value of θ that maximizes the log-likelihood function:

$$0 = (1/\theta) \cdot n - \sum_{i=1}^n \ln(x_i) \dots (I)$$

Solving for θ , we get:

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

Checking for maxima by taking a derivative w.r.t. θ again on (I):

$$\frac{\partial^2}{\partial \theta^2} \ln(L(\theta)) = -n / \theta^2 \dots (II)$$

Above value in eqn. (II) is negative, therefore indicating $\hat{\theta}_{MLE}$ is a maxima.

Thus, the maximum likelihood estimate of θ is the ratio of the number of observations n to the sum of 2 divided by each observation x_i .

Q2) To find the maximum a posteriori (MAP) estimator of θ , we need to find the value of θ that maximizes the posterior distribution $p(\theta | x_1, x_2, \dots, x_n)$. The posterior distribution is given by Bayes' theorem:

$$p(\theta | x_1, x_2, \dots, x_n) = (p(x_1, x_2, \dots, x_n | \theta) * p(\theta)) / p(x_1, x_2, \dots, x_n)$$

where $p(x_1, x_2, \dots, x_n | \theta)$ is the likelihood function and $p(\theta)$ is the prior distribution for θ .

The denominator $p(x_1, x_2, \dots, x_n)$ is a normalizing constant that ensures that the posterior distribution integrates to 1. Since we are only interested in the value of θ that maximizes the posterior distribution, we can ignore the normalizing constant and focus on the numerator of Bayes' theorem. The numerator is given by:

$$p(x_1, x_2, \dots, x_n | \theta) * p(\theta) = \theta^n \prod_{i=1}^n x_i^{-\theta-2} \lambda \theta^{-\lambda-2}$$

Taking the natural log of the numerator, we can simplify the problem to finding the value of θ that maximizes the log of the numerator:

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

Taking the derivative of the log of the numerator with respect to θ and setting it equal to zero, we can find the value of θ that maximizes the posterior distribution:

$$0 = (1/\theta) * n - \sum_{i=1}^n \ln(x_i) + (-\lambda - 2) / \theta \dots (I)$$

Solving for θ , we get:

$$\hat{\theta}_{\text{MAP}} = \frac{n - \lambda - 2}{\sum_{i=1}^n \ln(x_i)}$$

Checking for maxima by taking a derivative w.r.t. θ again on (I):

$$\frac{\partial^2}{\partial \theta^2} \ln(L(\theta)) = -(n - \lambda - 2) / \theta^2 \dots (II)$$

Above value in eqn. (II) is negative, therefore indicating $\hat{\theta}_{\text{MAP}}$ is a maxima.

Thus, the maximum a posteriori estimator of θ is the ratio of the difference between the number of observations n and the parameter governing the prior distribution λ to the sum of 2 divided by each observation x_i .

Q3) The five main steps to calculate the maximum likelihood estimate (MLE) given a distribution are:

1. Define the likelihood function. This is a function that represents the probability of a set of observations, given a set of parameters. The likelihood function is typically denoted as $L(\theta | x)$, where θ is the set of parameters and x is the set of observations.

$$L(\theta) = \prod f(x, \theta)$$

2. Define the log likelihood function. This is a function that represents the logarithm of the likelihood function. The log likelihood function is often easier to work with than the likelihood function because it is a monotonically increasing function, and it has the same maximum value as the likelihood function.

$$\ln(L(\theta)) = \ln(\prod f(x, \theta))$$

3. Specify the set of parameters that will be estimated. In most cases, the MLE will be a function of one or more unknown parameters, and these parameters will need to be estimated from the data. Choose an

optimization algorithm. There are many different optimization algorithms that can be used to find the maximum likelihood estimate, and the choice of algorithm will depend on the specific problem at hand. Some common optimization algorithms include gradient descent, Newton's method, and the Expectation-Maximization (EM) algorithm. In our case, we have differentiated to get the maxima of the function.

$$\frac{\partial}{\partial \theta} \ln(L(\theta))$$

4. Optimize the log likelihood function to find the maximum likelihood estimate. Once an optimization algorithm has been chosen, the log likelihood function can be optimized to find the set of parameters that maximize the likelihood of the observed data. This set of parameters is the maximum likelihood estimate.

$$step(3) = 0; find \hat{\theta}_{MLE}$$

5. Check if $\frac{\partial^2}{\partial \theta^2} \ln(L(\theta)) < 0$ [Maximized]

Once the maximum likelihood estimate has been found, it can be used to make predictions about future observations, or to perform hypothesis testing.

3: Naive Bayes Classifier

Q1) There are two types of independent parameters for a Naïve Bayes Classifier: Likelihood parameters & Prior parameters:

Likelihood parameters can be written as:

$$p(X = X_i | C = C_j)$$

where:

- X_i is the i^{th} feature, $i \in (1, 2, 3)$
- C_j is the j^{th} class, $j \in (1, 2)$

The total number of likelihood parameters will be $(3 - 1) * 2 = 4$. (The third likelihood parameter can be computed by using the other two likelihood parameters)

Prior parameters are one per class. The total number of prior parameters will be 1 from 2 classes since if we know prior of 1 class, we can compute the second prior.

Therefore, the total number of independent parameters will be **4 + 1 = 5**.

Q2) To predict the MLE for a class, we have the below relation:

$$c = \operatorname{argmax}_c \left(p(X = X_i | C = C_j) \times p(C = C_j) \right)$$

With the above relation, we will estimate the likelihood probabilities and the prior class probabilities below, using the data provided in the question:

For "Size" feature:

Size	Class = Yes	Class = No
Large	3	3
Small	1	3

Size	Class = Yes	Class = No
Large	0.75	0.5
Small	0.25	0.5

For "Color" feature:

Color	Class = Yes	Class = No	Color	Class = Yes	Class = No
Green	0	5	Green	0	0.833
Red	4	1	Red	1	0.167

For "Shape" feature:

Shape	Class = Yes	Class = No	Shape	Class = Yes	Class = No
Irregular	1	4	Irregular	0.25	0.667
Round	3	2	Round	0.75	0.333

For class priors:

Class	Prior	Class	Prior
Yes	4	Yes	0.4
No	6	No	0.6

Q3) Given $x = \{\text{Small, Red, Round}\}$,

$$p(y = \text{No}|x) = p(\text{Size} = \text{Small}|y = \text{No}) \times p(\text{Color} = \text{Red}|y = \text{No}) \times p(\text{Shape} = \text{Round}|y = \text{No}) \times p(y = \text{No})$$

$$p(y = \text{No}|x) = 0.5 \times 0.167 \times 0.333 \times 0.6$$

$$\mathbf{p(y = No|x) = 0.0166}$$

$$p(y = \text{Yes}|x) = p(\text{Size} = \text{Small}|y = \text{Yes}) \times p(\text{Color} = \text{Red}|y = \text{Yes}) \times p(\text{Shape} = \text{Round}|y = \text{Yes}) \times p(y = \text{Yes})$$

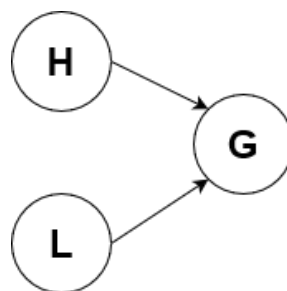
$$p(y = \text{Yes}|x) = 0.25 \times 1 \times 0.75 \times 0.4$$

$$\mathbf{p(y = Yes|x) = 0.075}$$

Since $p(y = \text{Yes}|x) > p(y = \text{No}|x)$, our Naïve Bayes classifier will predict $y = \text{Yes}$.

4: Probabilistic Graphical Models

Q1) In our case, the Bayesian Network is as follows:



The joint probability $p(G, H, L)$ is:

$$\mathbf{p(G, H, L) = p(H) \times p(L) \times p(G|H, L)}$$

Q2) From Q1. we know that:

$$p(G, H, L) = p(G|H, L) \cdot p(H) \cdot p(L)$$

On assuming no independence relations (given) we can deduce that:

$$p(H) = p(G, H, L) + p(\bar{G}, H, L) + p(G, H, \bar{L}) + p(\bar{G}, H, \bar{L})$$

$$\text{Hence, } p(H) = p_1 + p_2 + p_4 + p_6 \text{ (given) (1)}$$

$$p(L) = p(G, H, L) + p(\bar{G}, H, L) + p(G, \bar{H}, L) + p(\bar{G}, \bar{H}, L)$$

$$\text{Hence, } p(L) = p_1 + p_2 + p_3 + p_5 \text{ (given) (2)}$$

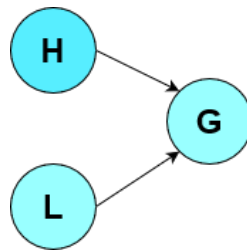
$$\text{Also, } p(G|H, L) = p(G, H, L) = p_1 \text{ (given) (3)}$$

Substituting (1), (2) and (3) in our Joint Probability equation we get:

$$p(G, H, L) = p_1 \cdot (p_1 + p_2 + p_4 + p_6) \cdot (p_1 + p_2 + p_3 + p_5)$$

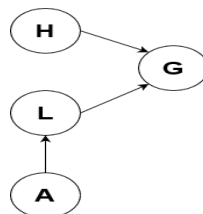
Hence, we can observe that in order to specify $p(G, H, L)$ we need 6 parameters ($p_1, p_2, p_3, p_4, p_5, p_6$)

Q3)



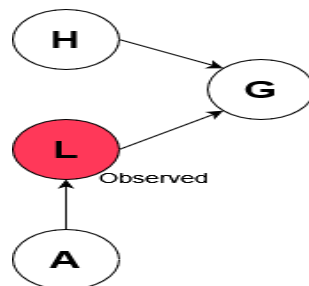
G is a **head-to-head** node and it is **observed**, therefore it is **unblocked**. Since there is an unblocked path between H and L, they are not independent.

Q4) Adding node for Lauren's assignments:



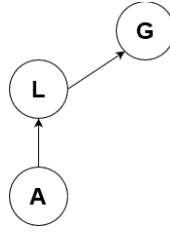
The joint probability $p(G, H, L, A)$ is

$$p(G, H, L, A) = p(H) \times p(A) \times p(L|A) \times p(G|H, L)$$



L is a **head-to-tail node** and it is **observed**, therefore it is **blocked**. A & G do not have any unblocked paths between them, and therefore are **independent**. Hence, if L is given, G and A are independent.

Q5) The new Bayesian Network after eliminating H is given as follows:



The joint probability $p(G, L, A)$ is:

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

On assuming no independence relations (given) we can deduce that:

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

$$\text{Hence, } p(A) = q_1 + q_2 + q_3 + q_5 \text{ (given) (1)}$$

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

$$\text{Hence, } p(L|A) = q_1 + q_2 \text{ (given) (2)}$$

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

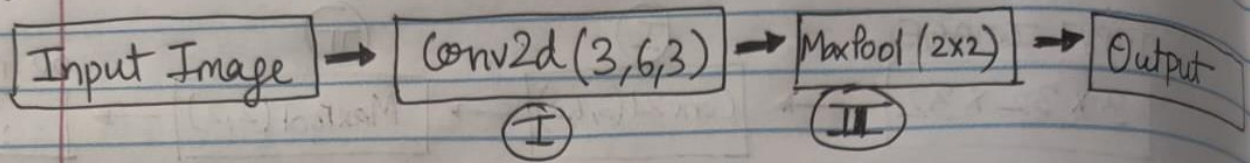
$$\text{Hence, } p(L|A) = q_1 + q_4 \text{ (given) (3)}$$

Substituting (1), (2) and (3) in our Joint Probability equation we get:

$$p(G, L, A) = (q_1 + q_4) \cdot (q_1 + q_2) \cdot (q_1 + q_2 + q_3 + q_5)$$

Hence, we can observe that in order to specify $p(G, L, A)$ we need 6 parameters: $(q_1, q_2, q_3, q_4, q_5)$

Q1)

5.1) Given Image Shape: $3 \times 32 \times 32$ 

① Size of Kernel in Conv2d: $3 \times 3 \Rightarrow K=3$

Stride: $1 \times 1 \Rightarrow S=1$ (Default)

Dimensions of Input Image: $D_1 \times H_1 \times W_1: 3 \times 32 \times 32$

Padding: $0 \Rightarrow P=0$ (Default)

Dimensions of Output Image: $D_2 \times H_2 \times W_2$

$$H_2 = (H_1 - K + 2P/S) + 1 = \frac{32-3}{1} + 1 = 30$$

$$W_2 = (W_1 - K + 2P/S) + 1 = \frac{32-3}{1} + 1 = 30$$

$$D_2 = \text{No. of filters} = 6$$

∴ Dimensions of Output Image: $(6 \times 30 \times 30)$ of (I)

② Size of Kernel in Max Pooling: $2 \times 2 \Rightarrow K=2$

Stride: $2 \times 2 \Rightarrow S=2$

Dimensions of Input Image: $D_2 \times H_2 \times W_2 = (6 \times 30 \times 30)$

Padding: $0 = P$

Dimensions of final Output Image: $D_3 \times H_3 \times W_3$

$$H_3 = (H_2 - K + 2P/S) + 1 = \frac{30-2}{2} + 1 = 15$$

$$W_3 = (W_2 - K + 2P/S) + 1 = \frac{30-2}{2} + 1 = 15$$

$$D_3 = 6$$

∴ Final size of the image will become $(6 \times 15 \times 15)$

Q2) Given Feature Map:

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

After applying a (2, 2) max pooling layer to the given feature map, we will get the following new feature map:

6	8
3	4

This is because the max pooling operation will take the maximum value from each 2x2 region of the input feature map and produce a new feature map with these maximum values.

If we change the max pooling layer to an average pooling layer, the resulting feature map would be:

3.25	5.25
2	2

This is because the average pooling operation will take the average value from each 2x2 region of the input feature map and produce a new feature map with these average values. Calculations:

1. $(1+1+5+6)/4 = 3.25$
2. $(2+4+7+8)/4 = 5.25$
3. $(3+2+1+2)/4 = 2$
4. $(1+0+3+4)/4 = 2$

Q3)

5.3) We are given,
 $x = -2, w = 5, u = -4$

Diagram illustrating the calculation of the output f from inputs x, w, u :

Inputs x and w are added to get $q = (x + w)$.
The result q is then multiplied by u to get the final output $f = (q * u)$.

Partial derivatives calculations:

$$\frac{\partial f}{\partial u} = \frac{\partial (q * u)}{\partial u} = q = (x + w) = 3$$
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \cdot \frac{\partial q}{\partial x} = \frac{\partial (q * u)}{\partial q} \cdot \frac{\partial (x + w)}{\partial x} = u = -4$$
$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial q} \cdot \frac{\partial q}{\partial w} = \frac{\partial (q * u)}{\partial q} \cdot \frac{\partial (x + w)}{\partial w} = u = -4$$

6:k-means

Q1)

There are a total of 6 points in the given set, and we want to choose 3 of them to be the initial centroids for the k-means algorithm. This can be done in $C_6^3 = 20$ ways. Therefore, there are 20 possible 3-starting configurations for the k-means algorithm.

Q2)

3-partition	Stable?	An example 3-starting configuration
{a, b}, {c, f}, {d, e}	Yes	{b, e, f}
{a, b, e}, {c, d}, {f}	No	none
{a, d}, {b, c}, {e, f}	Yes	{a, b, e}
{a, b, d}, {e}, {c, f}	No	none
{a, d, e}, {b, c}, {f}	Yes	{c, d, f}

7:Kernel Methods

Q1)

1) Given $x, x' \in \mathbb{R}^n$, a kernel $k: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is valid iff there exists a feature vector $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^m$ such that $k(x, x') = \phi(x)^T \phi(x')$.

1.1) To Prove: $k(x, x') = k_1(x, x') + k_2(x, x')$, where k_1 and k_2 are other valid kernels.

Since $k_1(x, x')$ & $k_2(x, x')$ are valid kernels:

$$k_1(x, x') = \phi^1(x)^T \phi^1(x')$$

$$k_2(x, x') = \phi^2(x)^T \phi^2(x')$$

$$\therefore k(x, x') = \phi^1(x)^T \phi^1(x') + \phi^2(x)^T \phi^2(x')$$

where, $\phi^1(x) = (\phi_1^1(x), \phi_2^1(x), \dots, \phi_{N_1}^1(x))$ } Feature Maps for $k_1 \in K$
 $\phi^2(x) = (\phi_1^2(x), \phi_2^2(x), \dots, \phi_{N_2}^2(x))$ }

We can therefore define a new kernel $\phi(x)$ where

$$\phi(x) = (\underbrace{\phi_1^1(x), \phi_2^1(x), \dots, \phi_{N_1}^1(x)}_{\text{Feature map for } k_1}, \underbrace{\phi_1^2(x), \phi_2^2(x), \dots, \phi_{N_2}^2(x)}_{\text{Feature map for } k_2})$$

This mapping satisfies

$$\phi(x) \cdot \phi(y) = \phi^1(x) \phi^1(y) + \phi^2(x) \phi^2(y)$$

$\therefore k(x, x')$ also is a valid kernel.

We can also prove it in Feature Vector Space Representation

$$\phi(x) = [\phi^1(x), \phi^2(x)]$$

$$k(x, x') = \langle \phi(x)^T, \phi(x') \rangle$$

$$= \langle [\phi^1(x)^T \phi^2(x)^T], [\phi^1(x') \phi^2(x')] \rangle$$

$$= \langle \phi^1(x)^T \phi^1(x') \rangle + \langle \phi^2(x)^T \phi^2(x') \rangle$$

$$\boxed{k(x, x') = k_1(x, x') + k_2(x, x')} \quad \because k(x, x') \text{ is a valid kernel.}$$

Q2)

7.2) To Prove: $k(x, x') = k_1(x, x') k_2(x, x')$ where k_1, k_2 are valid kernels.

Since $k_1(x, x')$ & $k_2(x, x')$ are valid kernels:

$$k_1(x, x') = \phi^1(x)^T \phi^1(x'); \quad k_2(x, x') = \phi^2(x)^T \phi^2(x')$$

$$\therefore k(x, x') = \phi^1(x)^T \phi^2(x)^T \phi^1(x') \phi^2(x')$$

This can be written as:

$$k(x, x') = \begin{bmatrix} \phi_1^1(x)^T \phi_1^2(x)^T \phi_1^1(x') \phi_1^2(x') & \dots & \phi_1^1(x)^T \phi_1^2(x)^T \phi_{N_2}^1(x') \phi_{N_2}^2(x') \\ \phi_2^1(x)^T \phi_2^2(x)^T \phi_1^1(x') \phi_1^2(x') & \dots & \phi_2^1(x)^T \phi_2^2(x)^T \phi_{N_2}^1(x') \phi_{N_2}^2(x') \\ \vdots & & \vdots \\ \phi_{N_1}^1(x)^T \phi_{N_1}^2(x)^T \phi_1^1(x') \phi_1^2(x') & \dots & \phi_{N_1}^1(x)^T \phi_{N_1}^2(x)^T \phi_{N_2}^1(x') \phi_{N_2}^2(x') \end{bmatrix}$$

↑
gram
matrix
representation

The above gram matrix representation is a Hadamard product (or element by element product) of K_1 & K_2 .

If K_1 is represented by a covariance matrix of (X_1, X_2, \dots, X_n)

If K_2 is represented by a covariance matrix of (Y_1, Y_2, \dots, Y_n)

then

$K = K_1 \odot K_2$ is represented simply by covariance matrix of $(X_1 Y_1, X_1 Y_2, \dots, X_n Y_n)$, implying that it is symmetric & positive definite.

K is also the Schur product.

We can also show this in Feature Vector Space Representation

$$\phi(x)_{ij} = \phi^1(x)_i \phi^2(x)_j$$

where $i=1, 2, \dots, N_1$ & $j=1, 2, \dots, N_2$

$$K(x, x') = \langle \phi(x)^T, \phi(x') \rangle$$

$$= \sum_{i=1}^{N_1} \sum_{j=1}^{N_2} \phi(x)^T_{ij} \phi(x')_{ij}$$

$$= \sum_{i=1}^{N_1} \phi^1(x)^T_i \phi^1(x')_i \sum_{j=1}^{N_2} \phi^2(x)^T_j \phi^2(x')_j$$

$$K(x, x') = K_1(x, x') K_2(x, x')$$

∴ $K(x, x')$ is a valid kernel.

Q3)

1.3) To Prove: $k(x, x') = f(x) k_1(x, x') f(x')$, where $f: \mathbb{R}^n \rightarrow \mathbb{R}$

Since $k(x, x')$ is the product of $f(x)$, $k_1(x, x')$, and $f(x')$, we can write $k(x, x')$ as:

$$k(x, x') = f(x) \underbrace{(\phi'(x)^T \phi'(x'))}_{k_1 \text{ is a valid kernel}} f(x')$$

$$= \underbrace{f(x) \phi'(x)^T}_{\phi(x)^T} \underbrace{\phi'(x') f(x')}_{\phi(x')}$$

Expression (I) is of the form $(\phi(x)^T * \phi(x'))$ where $\phi(x)$ is the product of the scalar function f and the feature vector of x corresponding to k_1 .

∴ We have shown that there exists a feature vector $\phi(x)$ such that $k(x, x') = \phi(x)^T \phi(x')$. ϕ is just another kernel transformation in that domain.

$$\boxed{\phi(x) = (f(x) \phi'_1(x), f(x) \phi'_2(x), \dots, f(x) \phi'_N(x))}$$

∴ $k(x, x')$ is a valid kernel.

7.4) To Prove: $K(x, x') = \exp \{ k_1(x, x') \}$, where k_1 is another valid kernel

We will first show

$$K(x, x') = \alpha K_1(x, x') + \beta K_2(x, x') \quad \dots \textcircled{I}$$

where K_1 & K_2 are valid kernels

① Since $\alpha K_1(x, x') = \langle \sqrt{\alpha} \phi^1(x), \sqrt{\alpha} \phi^1(x') \rangle \dots \textcircled{A}$

$$\beta K_2(x, x') = \langle \sqrt{\beta} \phi^2(x), \sqrt{\beta} \phi^2(x') \rangle \dots \textcircled{B}$$

Then, $K(x, x') = \langle \sqrt{\alpha} \phi^1(x), \sqrt{\alpha} \phi^1(x') \rangle + \langle \sqrt{\beta} \phi^2(x), \sqrt{\beta} \phi^2(x') \rangle$
using ① & ②

$$K(x, x') = \langle [\sqrt{\alpha} \phi^1(x) \sqrt{\beta} \phi^2(x)], [\sqrt{\alpha} \phi^1(x') \sqrt{\beta} \phi^2(x')] \rangle$$

∴ K can be expressed as an inner product and is a valid kernel.

② We also know from the proof in 7.2 that
 $K(x, x') = K_1(x, x') K_2(x, x') \dots \textcircled{II}$

③ Now, we will show

$$K(x, x') = f(K_1(x, x')) \quad \text{where } f \text{ is a polynomial with valid kernel +ve co-efficients}$$

Since each polynomial term is a product of kernels with a +ve co-efficient using ① & ②, we can deduce that

$$K(x, x') = f(K_1(x, x')) \text{ is a valid kernel } \textcircled{III}$$

Now to prove that $K(x, x') = \exp \{ k_1(x, x') \}$ we can write

$$\exp(x) = \lim_{i \rightarrow \infty} \left(1 + x + \dots + \frac{x^i}{i!} \right) \quad (\text{Taylor Series Expansion})$$

Now using ③, we can say

$$K(x, x') = \lim_{i \rightarrow \infty} K_i(x, x')$$

Since $K(x, x')$ is a sum of valid kernels ∴ it is a valid kernel.

Q5)

7.5) To Prove: $k(x, x') = \exp \left\{ -\frac{\|x - x'\|^2}{2\sigma^2} \right\}$

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

$$\therefore k(x, x') = \exp\left(\frac{-x^T x}{2\sigma^2}\right) \exp\left(\frac{x^T x'}{\sigma^2}\right) \exp\left(\frac{-(x')^T x'}{2\sigma^2}\right) \dots \textcircled{A}$$

We have proven that below kernels are valid:

(I) $k(x, x') = x^T x' \Rightarrow$ Linear Kernel \Rightarrow Rename it to $k_1(x, x')$

(II) $k(x, x') = \exp(k_1(x, x')) \Rightarrow$ From 7.4

(III) \therefore Eqn (A) $\Rightarrow k(x, x') = f(x) k_1(x, x') f(x')$ is valid kernel

New Kernel $\Rightarrow k(x, x') = \exp \left\{ -\frac{1}{2\sigma^2} (k(x, x') + k(x', x') - 2k(x, x')) \right\}$

Q2)

8.2) Given: $p(x) = \sum_{k=1}^K \pi_k p(x/k)$

Covariances of k_1, k_2, \dots, k_n distributions are $\Sigma_1, \Sigma_2, \dots, \Sigma_n$.

$$E[x] = \sum_{k=1}^K \pi_k \mu_k \quad \text{proved above in 8.1}$$

Mean of k_1, k_2, \dots, k_n distributions are $\mu_1, \mu_2, \dots, \mu_n$.

We can write: $\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$

Conventions: Let there be a random variable θ with possible values $1, \dots, K$ and probabilities $\pi_1, \pi_2, \dots, \pi_K$.
mixing co-efficient

$$E[X|\theta=k] = \mu_k \quad \dots \textcircled{I}$$

$$\text{Cov}(X|\theta=k) = \Sigma_k \quad \dots \textcircled{II}$$

$$(\Sigma_k)_{mn} = \text{Cov}(X_m, X_n | \theta=k) = E(X_m X_n | \theta=k) - E(X_m | \theta=k) E(X_n | \theta=k) \quad \textcircled{A}$$

We have already shown $E(X) = \sum_{k=1}^K \mu_k \pi_k \quad \dots \textcircled{B}$

$$E(X_m X_n) = \sum_{k=1}^K E(X_m X_n | \theta=k) \pi_k$$

$$\textcircled{C} \quad E(X_m X_n) = \sum_{k=1}^K \left((\Sigma_k)_{mn} + (\mu_k)_m (\mu_k)_n \right) \pi_k \quad \dots \text{using } \textcircled{A} \text{ \& } \textcircled{B}$$

Continuation:

Using (C) & (B) \rightarrow

$$\text{Cov}(X_m X_n) = E(X_m X_n) - E(X_m) E(X_n)$$

$$= \sum_{k=1}^K ((\Sigma_k)_{mn} + (\mu_k)_m (\mu_k)_n) - \sum_{k=1}^K (\mu_k)_m \pi_k \sum_{k=1}^K (\mu_k)_n \pi_k$$

This can be expressed as

$$\text{Cov}(X_m X_n) = E[\text{Cov}(X_m X_n | \theta)] + \text{Cov}(E(X_m | \theta), E(X_n | \theta))$$

"Law of Total Covariance"

$$\text{Cov}[X] = \sum_{k=1}^K \Sigma_k \pi_k + \sum_{k=1}^K \mu_k^T \mu_k \pi_k - \underbrace{\mu^T \mu}_{E[X] E[X]^T}$$

$$\boxed{\text{Cov}[X] = \sum_{k=1}^K \pi_k (\Sigma_k + \mu_k^T \mu_k) - E[X] E[X]^T}$$

Hence Proved

9: Approximate Inference

Q1)

9.1) We are given a fully Bayesian Model

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

$$Z = \{z_1, z_2, \dots, z_N\} \text{ latent variables \& parameters}$$

Model specifies joint distribution $p(X, Z)$

Using Bayes' Theorem:

$$p(X) = \frac{p(X, Z)}{p(Z|X)} = \frac{\frac{p(X, Z)}{q(Z)}}{\frac{p(Z|X)}{q(Z)}}$$

where $q(Z)$ is the distribution of the latent variables

Taking \ln on both sides:

$$\ln(p(X)) = \ln\left(\frac{p(X, Z)}{p(Z|X)}\right) = \ln\left(\frac{p(X, Z)}{q(Z)}\right) - \ln\left(\frac{p(Z|X)}{q(Z)}\right)$$

Multiplying by $q(Z)$ on both sides:

$$q(Z) \ln(p(X)) = q(Z) \ln\left(\frac{p(X, Z)}{q(Z)}\right) - q(Z) \ln\left(\frac{p(Z|X)}{q(Z)}\right)$$

Integrating w.r.t. Z on both sides:

$$\int q(Z) \ln(p(X)) dZ = \int q(Z) \ln\left(\frac{p(X, Z)}{q(Z)}\right) dZ - \int q(Z) \ln\left(\frac{p(Z|X)}{q(Z)}\right) dZ$$

\downarrow \downarrow \downarrow

$\ln(p(X)) \int q(Z) dZ$ $L(q)$ given $KL(q||p)$ given

Since $q(Z)$ is a distribution, $\int q(Z) dZ = 1$
so we get,

$$\boxed{\ln(p(X)) = L(q) + KL(q||p)}$$

10: Hidden Markov Models

Q1) To Prove: If π or A are initially set to zero, then they will remain zero in all subsequent updates of the EM algorithm.

EM algorithm for a Hidden Markov Model:

E-step: Computation of the expectation of the complete-data log likelihood for model parameters

M-step: Updating model parameters to maximize expectation

For E-step: Since the initial values of the model parameters include elements that are set to zero, these zero elements will not contribute to the expected value of the complete-data log likelihood. Therefore, the updated model parameters in the M-step will not be affected by these zero elements.

For M-step: Since the initial values of the model parameters include elements that are set to zero, these zero elements will not affect the maximization of the expected value of the complete-data log likelihood. Therefore, the updated model parameters will not be affected by these zero elements, and they will remain zero in all subsequent updates of the EM algorithm.

π and A updated based on the following update rules:

$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})}$$
$$A_{jk} = \frac{\sum_{t=1}^T \xi(z_{t-1}, j, z_{tk})}{\sum_{m=1}^K \sum_{t=2}^T \xi(z_{t-1}, j, z_{tm})}$$
$$\xi(z_{t-1}, z_t) = \frac{\alpha(z_{t-1}) P(x_t | z_t) P(z_t | z_{t-1}) \beta(z_t)}{P(x)}$$

If $A_{jk} = 0$, then $\xi(z_{t-1}, z_t) = 0$

∴ All subsequent updates in the EM algorithm are 0.

Therefore, if any elements of the parameter π or A for a Hidden Markov Model are initially set to zero, then those elements will remain zero in all subsequent updates of the EM algorithm.

Q2) A Markov model is a mathematical model used to describe the behavior of a system where the future state of the system is dependent only on its current state, not on its past history. In this cake factory example, the state of the system is determined by the positions of the Chocolate and Caramel levers, and the future state of the system is determined by the actions of the controller staff in the control room. The possible states of the system are as follows:

- Chocolate lever on, Caramel lever on (Chocolate+Caramel cake) **(ChCa)**
- Chocolate lever on, Caramel lever off (Chocolate cake) **(Ch)**
- Chocolate lever off, Caramel lever on (Caramel cake) **(Ca)**
- Chocolate lever off, Caramel lever off (Plain cake) **(P)**

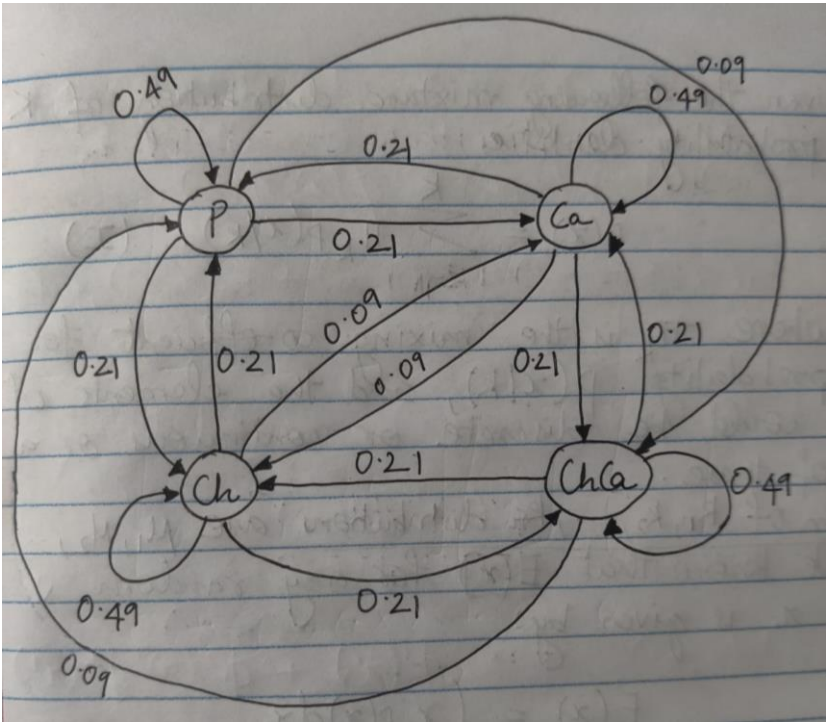
The prior state distribution is the initial probability distribution over the possible states of the system. In this case, since the Chocolate and Caramel levers are independently switched on or off at the beginning of the day with a 50% chance each, the prior state distribution is as follows:

- Chocolate lever on, Caramel lever on: 0.25
- Chocolate lever on, Caramel lever off: 0.25
- Chocolate lever off, Caramel lever on: 0.25
- Chocolate lever off, Caramel lever off: 0.25

This transition matrix indicates the probabilities of transitioning from one state to another when the controller staff in the control room switches the positions of the Chocolate and Caramel levers. The state transition probabilities are ((the probability of both levers being switched is 0.09, the probability of one being switched is $0.7 \times 0.3 = 0.21$, and the probability of neither being switched is 0.49).

	P	Ch	Ca	ChCa
P	0.49	0.21	0.21	0.09
Ch	0.21	0.49	0.09	0.21
Ca	0.21	0.09	0.49	0.21
ChCa	0.09	0.21	0.21	0.49

The Markov model for the cake factory can be represented as a directed graph, where the nodes represent the possible states of the system and the edges represent the transitions between states. The prior state distribution can be represented as the initial probabilities of being in each state, and the transition matrix can be represented as the probabilities of transitioning from one state to another along the edges of the graph. For example, the Markov model for the cake factory can be represented as follows:



Q3) The probability that the machine will produce the sequence {Plain, Chocolate, Chocolate, Chocolate+Caramel} in this order depends on the initial state of the system and the specific sequence of transitions between the states.

Assuming the initial state is randomly selected from the prior state distribution, the probability of producing the sequence {Plain, Chocolate, Chocolate, Chocolate+Caramel} in this order is as follows:

1. The initial state is Plain, with a probability of 0.25.
2. The next state is Chocolate, with a probability of 0.21 (transitioning from Plain to Chocolate with a probability of 0.21 according to the transition matrix).
3. The next state is Chocolate, with a probability of 0.49 (transitioning from Chocolate to Chocolate with a probability of 0.49 according to the transition matrix).
4. The next state is Chocolate+Caramel, with a probability of 0.21 (transitioning from Chocolate to Chocolate+Caramel with a probability of 0.21 according to the transition matrix).

The overall probability of producing the sequence {Plain, Chocolate, Chocolate, Chocolate+Caramel} is **$P(P, Ch, Ch, ChCa) = P(P) * P(Ch|P) * P(Ch|Ch) * P(ChCa|Ch)$** in this order is the product of the individual probabilities:

which is equal to **$0.25 * 0.21 * 0.49 * 0.21 = 0.0054$** .

Note that this is just one possible sequence of states that can result in the production of the sequence {Plain, Chocolate, Chocolate, Chocolate+Caramel} of cakes. There may be other sequences of states that result in the same sequence of cakes with different probabilities.