

1. Find a SVM classifier for the given data and classify the point  $(4, -4)$ .

$x_1$	$x_2$	class label
6	-2	+
12	2	+
6	2	+
12	-2	+
0	-2	-
-2	0	-
2	0	-
0	2	-

The support vectors can be

$$s_1 = \begin{bmatrix} 6 \\ -2 \end{bmatrix}$$

$$s_2 = \begin{bmatrix} 6 \\ 2 \end{bmatrix}$$

$$s_3 = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

$$y_1 = +1$$

$$y_2 = +1$$

$$y_3 = -1$$

$$\sum_{i=1}^3 \alpha_i y_i = 0$$

$$\alpha_1(+1) + \alpha_2(+1) + \alpha_3(-1) = 0$$

$$\alpha_3 = \alpha_1 + \alpha_2$$

$$w x + b = -1 \Rightarrow w s_3 + b = -1$$

$$2w_1 + b = -1 \quad \text{--- (1)}$$

$$w x + b = +1 \Rightarrow w s_1 + b = 1$$

$$6w_1 - 2w_2 + b = 1 \quad \text{--- (2)}$$

$$ws_2 + b = 1$$

$$6w_1 + 2w_2 + b = 1 \quad \text{--- (3)}$$

$$(2) + (3) \Rightarrow 12w_1 + 2b = 2 \quad \text{--- (4)}$$

$$(4) - 2(1) \Rightarrow 8w_1 = 4$$

$$w_1 = \frac{1}{2} \Rightarrow b = -2, w_2 = 0$$

$$w = \begin{bmatrix} 1/2 \\ 0 \end{bmatrix} \quad b = -2$$

$$w = \begin{bmatrix} 1/2 \\ 0 \end{bmatrix}$$

$$b = -2$$

give  $x = [4, -4]$

$$\begin{aligned} w \cdot x + b &= 1/2(4) - 4(0) - 2 \\ &= 2 - 2 \\ &= 0 \end{aligned}$$

$$w \cdot x + b = 0$$

add an isolated term  $\rightarrow$

$$\begin{bmatrix} 1/2 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 4 \\ -4 \end{bmatrix} + b = \begin{bmatrix} 1/2 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 4 \\ -4 \end{bmatrix} + 12$$

$$1/2 \cdot 4 + 0 + b = 1/2 \cdot 4 + 12$$

$$2 + b = 12$$

$$b = 12 - 2$$

$$b = 10$$

$$b = 10 \rightarrow \text{add } 10 \text{ to both sides}$$

$$10 + 10 = 10 + 10$$

$$10 + 10 = 10 + 10 \rightarrow \text{cancel}$$

$$10 + 10 = 10 + 10$$

$$10 + 10 = 10$$

$$10 + 10 = 10 + 10$$

$$10 + 10 = 10 + 10 \rightarrow \text{cancel}$$

$$10 + 10 = 10 + 10$$

$$10 + 10 = 10 + 10 \rightarrow \text{cancel}$$

$$10 + 10 = 10 + 10$$

$$10 + 10 = 10 + 10$$

$$10 + 10 = 10 + 10$$

2 consider the following 2D data points: Find out the points that belongs to each cluster by assuming  $k=3$  and  $A_4, A_8$  and  $A_{12}$  are the cluster centers of  $C_1, C_2$  and  $C_3$  respectively. Perform 2 iterations.

Point	co-ordinates
A1	(2,10)
A2	(2,6)
A3	(11,11)
A4	(6,9)
A5	(6,4)
A6	(1,2)
A7	(5,10)
A8	(4,9)
A9	(10,12)
A10	(7,5)
A11	(9,11)
A12	(4,6)
A13	(3,10)
A14	(3,8)
A15	(6,11)

#### Iteration 1

Data points	Distance to $C_1$	Distance to $C_2$	Distance to $C_3$
A1	6.08	3.16	4.47
A2	4.12	2.24	2.83
A3	8.06	2.24	5.10
A4	0	3.61	5.10
A5	5.66	2.24	3.61
A6	7.07	5.83	2.24
A7	3.16	0.41	5.83

A8	2.24	0	4.12	CA
A9	5.39	6.71	3.16	SA
A10	5.83	4.47	2.24	FA
A11	4.12	0	3.16	FA
A12	2.24	3.61	0	FA
A13	4.12	0.41	4.47	FA
A14	4.47	1.41	5.10	SA
A15	5.37	0	3.16	FA

Data point	Assigned cluster	cluster	New center
A1	C3	C1	(5.5, 9.5)
A2	C3	C1	(6.6, 6.2)
A3	C2	C2	(4, 6)
A4	C1	C3	
A5	C2	C2	
A6	C3	C3	
A7	C1	C2	
A8	C2	C2	
A9	C2	C2	
A10	C2	C2	
A11	C2	C2	
A12	C3	C3	
A13	C3	C3	
A14	C3	C3	
A15	C2	C2	

## Iteration 2

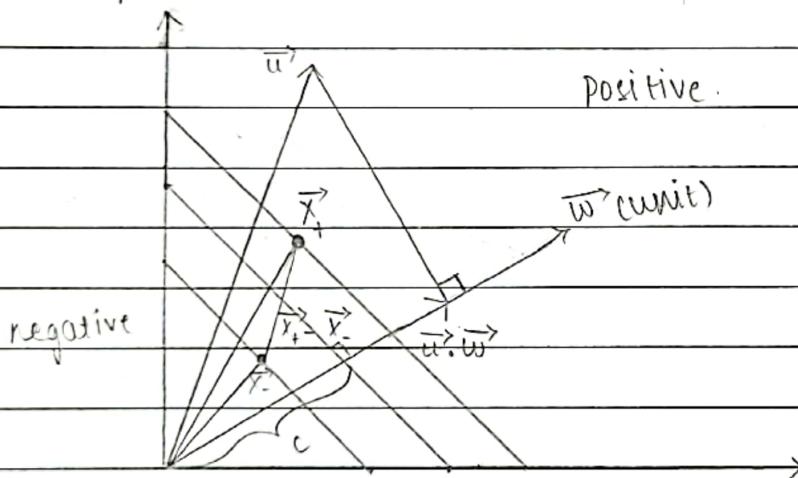
Data Point	Distance to C1	Distance to C2	Distance to C3
A1	6.92	3.20	3.16

A2	4.58	2.32	2.2	NA
A3	9.11	2.32	4.47	NA
A4	2.83	3.88	4.47	NA
A5	4.74	2.32	3.88	NA
A6	7.81	6.10	1.41	NA
A7	2	0.89	5.83	NA
A8	0.89	0	3.16	NA
A9	4.58	7.79	2.32	NA
A10	5.39	4.47	1.41	
A11	2.83	0	2.32	
A12	0.89	3.16	0	NA
A13	3.16	0.89	2.32	NA
A14	3.6	0	0	NA

A15

Data point	Assigned cluster	cluster	center
A1	C3	C1	(11, 11)
A2	C3	C1	(11, 11)
A3	C2	C2	(6, 6)
A4	C1	C3	(3, 3)
A5	C2	C2	(6, 6)
A6	C3	C3	(6, 6)
A7	C1	C1	(11, 11)
A8	C2	C2	(6, 6)
A9	C3	C3	(6, 6)
A10	C3	C3	(6, 6)
A11	C2	C2	(6, 6)
A12	C3	C3	(6, 6)
A13	C2	C2	(6, 6)
A14	C1	C1	(11, 11)

3. The SVM hard margin classifier aims to find the optimal hyperplane that maximizes the margin between the two classes while ensuring that all training points are correctly classified. Derive the following:
- The optimization problem that needs to be solved for finding the optimal hyperplane.
  - The lagrangian for the optimization problem.
  - The dual form of the optimization problem.
  - The equations for  $w$  and  $b$  in terms of the Lagrange multipliers.



(a) In other words,

$$\vec{w} \cdot \vec{u} \geq c \quad \text{--- positive}$$

$$\vec{w} \cdot \vec{u} \leq c \quad \text{--- negative}$$

$$\vec{w} \cdot \vec{u} - c \geq 0 \quad \text{--- (1)}$$

$$\vec{w} \cdot \vec{u} - c \leq 0 \quad \text{--- (2)}$$

Let  $\vec{u} = \vec{x}_+$  for positive sample and  $\vec{x}_-$  for negative sample

At the boundary, let the  $w$  constants be  $c_1$  and  $c_2$

$$\vec{w} \cdot \vec{x}_+ - c_1 = 0 \quad \text{--- (3)} \qquad \vec{w} \cdot \vec{x}_- - c_2 = 0 \quad \text{--- (4)}$$

$$\text{Let } A = -\frac{(c_2 + c_1)}{2} \quad \text{and} \quad B = -\frac{(c_2 - c_1)}{2}$$

$$\text{Then } -c_1 = A - B \quad \text{and} \quad c_2 = A + B$$

Replacing in (3) and (4)

$$\vec{w} \cdot \vec{x}_+ + A - B = 0 \quad \text{--- (5)}$$

$$\vec{w} \cdot \vec{x}_- + A + B = 0 \quad \text{--- (6)}$$

Dividing (5) and (6) by B

$$\frac{\vec{w} \cdot \vec{x}_+ + A}{B} - 1 = 0 \quad \text{--- (7)}$$

$$\frac{\vec{w} \cdot \vec{x}_- + A}{B} + 1 = 0 \quad \text{--- (8)}$$

Redefining  $\frac{\vec{w}}{B} = \vec{w}$  and  $A = b$  in (7) and (8)

$$\vec{w} \cdot \vec{x}_+ + b = 1 \quad \text{--- (9)}$$

$$\vec{w} \cdot \vec{x}_- + b = -1 \quad \text{--- (10)}$$

converting to inequality

$$\vec{w} \cdot \vec{x}_+ + b \geq 1 \quad \text{--- (11)}$$

$$\vec{w} \cdot \vec{x}_- + b \leq -1 \quad \text{--- (12)}$$

The predicted values of y

$$y_i = 1 \quad : +ve$$

$$y_i = -1 \quad : -ve$$

Multiplying (11) and (12) with  $y_i$

$$y_i (\vec{w} \cdot \vec{x}_+ + b) - 1 \geq 0 \quad \text{--- (13)}$$

$$y_i (\vec{w} \cdot \vec{x}_- + b) - 1 \leq 0 \quad \text{--- (14)}$$

Generalizing for all  $\vec{x}$ 's

$$y_i (\vec{w} \cdot \vec{x} + b) - 1 \geq 0 \quad \text{--- (15)}$$

### (b) Multiple constraints

$$L(x, \lambda) = f(x) + \lambda_1 g_1(x) + \lambda_2 g_2(x) + \dots + \lambda_K g_K(x)$$

Objective  $f(w) = \frac{1}{2} \|w\|^2$  to minimize

Subject to m constraints  $g_i(w, b) = y_i (w^T x_i + b) - 1$

$$L = f(w) - \sum_{i=1}^m y_i g_i(w, b)$$

$$L = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y_i (w^T x_i + b) - 1] \quad (16)$$

(d) Derivative wrt  $w$

$$\frac{\partial L}{\partial w} = \vec{w} - \sum \alpha_i y_i \vec{x}_i = 0$$

$$\vec{w} = \sum_{i=1}^m \alpha_i y_i \vec{x}_i \quad (17)$$

Derivative wrt  $b$

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^m \alpha_i y_i = 0$$

$$\sum_{i=1}^m \alpha_i y_i = 0 \quad (18)$$

(c) Replacing (17) and (18) in (16)

$$L = \frac{1}{2} \left( \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j \right) - \left( \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j - b \sum_{i=1}^m \alpha_i y_i \right)$$

$$+ \sum_{i=1}^m \alpha_i$$

$$L(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \left( \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j \right)$$

$$\max_w L(w) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \left( \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j \right)$$

4. suppose you have a 40% chance of it raining (event A) and there is a 60% chance that it will be cloudy (event B) on any given day. You also know that when it rains, there's a 45% chance of it being cloudy.

$$\text{Bayes' theorem: } P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Here event A is raining.

Event B is being cloudy.

$P(A)$  - probability of rain 40% or 0.4

$P(B)$  - probability of being cloudy 60% or 0.6

$P(B|A)$  - probability of being cloudy given that it is raining 45% or

$$P(A) \cdot P(B|A) = 0.4 \cdot 0.45 = 0.18$$

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} = \frac{0.18}{0.6} = 0.3$$

$$= 0.45 \times 0.4$$

$$0.6 \cdot (A|I=X) \cdot (A|I=Y) = (A|I=X) \cdot (A|I=Y)$$

$$P(A|B) = 0.3$$

∴ probability that it is raining given that it is cloudy is 30%.

$$(A|I=X) \cdot (A|I=Y) = (A|I=X) \cdot (A|I=Y)$$

$$P(A|I=X) \cdot P(A|I=Y)$$

$$P(A|I=X) + P(A|I=Y)$$

$$P(A|I=X) + P(A|I=Y)$$

$$P(A|I=X) + P(A|I=Y)$$

5. suppose you have two classes, class A and class B, and two features (X and Y). For class A, the conditional probabilities of the features are as follows:

$$P(X=1|A) = 0.2 \quad P(Y=1|A) = 0.3$$

For class B, the conditional probabilities of the features are

$$P(X=1|B) = 0.8 \quad P(Y=1|B) = 0.7$$

The prior probabilities of the classes are:

$$P(A) = 0.3$$

$$P(B) = 0.7$$

You're given a new instances with features  $X=1$  and  $Y=1$ . use the Baye's optimal classifier to determine which class is most probable for this instance.

$$\text{for } X=1, Y=1 \text{ using point } 1 \text{ in notes} \rightarrow P(A|X=1, Y=1)$$

$$P(A|X=1, Y=1) = \frac{P(X=1, Y=1|A) \cdot P(A)}{P(X=1, Y=1)}$$

$$P(B|X=1, Y=1) = \frac{P(X=1, Y=1|B) \cdot P(B)}{P(X=1, Y=1)}$$

- For Class A

$$P(X=1, Y=1|A) = P(X=1|A) \cdot P(Y=1|A) \\ = 0.2 \times 0.3$$

$$P(X=1, Y=1|A) = 0.06$$

$$P(A|X=1, Y=1) = 0.06 \times 0.3$$

$$P(A|X=1, Y=1) = 0.018$$

- For Class B

$$P(A|X=1, Y=1|B) = P(X=1|B) \cdot P(Y=1|B) \\ = 0.8 \times 0.7$$

$$P(X=1, Y=1|B) = 0.56$$

$$P(B|X=1, Y=1) = 0.56 \times 0.7$$

$$P(B|X=1, Y=1) = 0.392$$

Normalizing the probabilities:

$$P(A|X=1, Y=1) = \frac{0.018}{0.018 + 0.392} = \frac{0.018}{0.41} = 0.0439$$

$$P(B|X=1, Y=1) = \frac{0.392}{0.018 + 0.392} = \frac{0.392}{0.41} = 0.956$$

According to the Bayes optimal classifier the new instance with features  $X=1$  and  $Y=1$  is most probable to belong to class B

6. consider the following data:

Text

Class

It was a bad movie

-

the movie had no plot

-

It was a great movie

+

the movie had a good plot

+

the movie was a boring movie

-

classify the text "the movie had a boring plot"

Note: consider stop words for this.

The text "the movie had a boring plot" can be classified as negative.

(i) Remove stop words: stop words are words that are common in a language and do not add much meaning to a text. In this case the stop words are 'the', 'a', and 'had'. After removing the stop words, the text is "movie boring plot".

(ii) The word boring has a negative sentiment, so sentiment of the text is negative.

∴ the text "The movie had a boring plot" can be classified as negative.

7. given the following states and their transition and emission probabilities, calculate the probability of the output sequence

$Q_1 \rightarrow Q_3 \rightarrow Q_2$  using the forward algorithm by using the given HMM parameters

$$\pi_A = 0.3, \pi_B = 0.2, \pi_C = 0.5$$

transition and emission table:

	A	B	C	01	02	03	
A	0.1	0.5	0.4	0.2	0.2	0.6	
B	0.7	0.2	0.1	0.4	0.4	0.2	
C	0.1	0.1	0.8	0.5	0.4	0.1	

where  $01, 02, 03$  are the emission states.

(i) Initialize the forward probabilities:

$$\alpha_1(A) = \pi_A \times p(01|A)$$

$$= 0.3 \times 0.2$$

$$\alpha_1(A) = 0.06$$

$$\alpha_1(B) = \pi_B \times p(01|B)$$

$$= 0.2 \times 0.4$$

$$\alpha_1(B) = 0.08$$

$$\alpha_1(C) = \pi_C \times p(01|C)$$

$$= 0.5 \times 0.5$$

$$\alpha_1(C) = 0.25$$

(ii) calculate the forward probabilities for time step  $t=2$

$$\alpha_2(A) = \alpha_1(A) \times p(B|A) \times p(02|B)$$

$$= 0.06 \times 0.5 \times 0.4$$

$$\alpha_2(A) = 0.012$$

$$\alpha_2(B) = \alpha_1(A) \times p(C|B) \times p(02|C)$$

$$= 0.08 \times 0.1 \times 0.4$$

$$\alpha_2(B) = 0.0032$$

$$\alpha_2(C) = \alpha_1(C) \times P(A|C) \times P(O_2|A)$$

$$= 0.25 \times 0.8 \times 0.2 = 0.16$$

$$\alpha_2(C) = 0.04$$

(3) calculate the forward probabilities for time step  $t=3$ :

$$\alpha_3(A) = \alpha_2(A) \times P(C|A) \times P(O_3|C)$$

$$= 0.012 \times 0.4 \times 0.6$$

$$\alpha_3(A) = 0.00288$$

$$\alpha_3(B) = \alpha_2(B) \times P(A|B) \times P(O_3|A)$$

$$= 0.0032 \times 0.7 \times 0.2$$

$$\alpha_3(B) = 0.000448$$

$$\alpha_3(C) = \alpha_2(C) \times P(B|C) \times P(O_3|B)$$

$$= 0.04 \times 0.1 \times 0.2$$

$$\alpha_3(C) = 0.0008$$

(4) calculate the probability of the output sequence

$$P(O_1 \rightarrow O_3 \rightarrow O_2) = \alpha_3(A) + \alpha_3(B) + \alpha_3(C) =$$

$$= 0.00288 + 0.000448 + 0.0008 = 0.004128$$

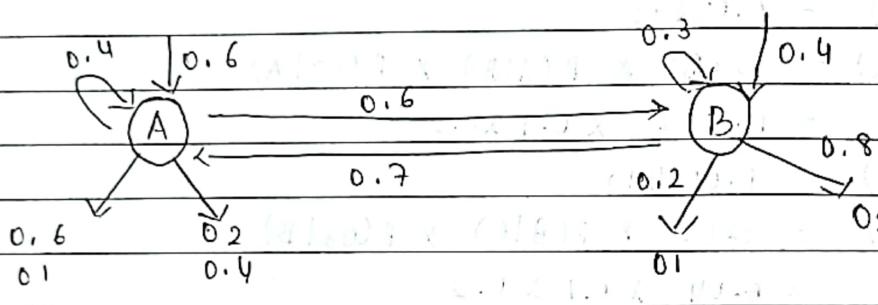
$$P(O_1 \rightarrow O_3 \rightarrow O_2) = 0.004128$$

$\therefore$  The probability of the output sequence  $O_1 \rightarrow O_3 \rightarrow O_2$  is 0.004128.

8. Given the following HMM, use Viterbi's algorithm to compute the most probable final hidden state given the following output sequence is generated:

$0_2 \rightarrow 0_1 \rightarrow 0_1$

HMM:



(1) Initialize the Viterbi probabilities

$$\delta_1(A) = \pi_A \times P(0_2|A)$$

$$= 0.3 \times 0.2$$

$$\delta_1(A) = 0.06$$

$$\delta_1(B) = \pi_B \times P(0_2|B)$$

$$= 0.2 \times 0.4$$

$$\delta_1(B) = 0.08$$

$$\delta_1(C) = \pi_C \times P(0_2|C)$$

$$= 0.5 \times 0.1$$

$$\delta_1(C) = 0.05$$

(2) calculate the viterbi probabilities and back pointers for time

Step  $t = 2$

$$\delta_2(A) = \max \{ \delta_1(A) \times P(A|A) \times P(0_1|A), \delta_1(B) \times P(B|A) \times P(0_1|B), \delta_1(C) \times P(C|A) \times P(0_1|C) \}$$

$$= \max \{ 0.06 \times 0.1 \times 0.2, 0.08 \times 0.7 \times 0.4, 0.05 \times 0.8 \times 2.4 \}$$

$$\delta_2(A) = 0.0224$$

$$\delta_2(B) = \max \{ \delta_1(A) \times P(A|B) \times P(0_1|A), \delta_1(B) \times P(B|B) \times P(0_1|B) \}$$

$$\begin{aligned} & s_1(c) \times p(c|b) \times p(o_1|c) \\ &= \max \{0.06 \times 0.5 \times 0.2, 0.8 \times 0.2 \times 0.4, 0.05 \times 0.8 \times 0.4\} \\ &= 0.0192 \end{aligned}$$

$$\begin{aligned} s_2(c) &= \max \{s_1(a) \times p(a|c) \times p(o_1|a), s_1(b) \times p(b|c) \times p(o_1|b), \\ &\quad s_1(c) \times p(c|c) \times p(o_1|c)\} \\ &= \max \{0.06 \times 0.8 \times 0.2, 0.08 \times 0.1 \times 0.4, 0.05 \times 0.8 \times 0.5\} \\ &= 0.02 \end{aligned}$$

(3) calculate the viterbi probabilities and back pointers for time

Step:  $t = 3$

$$\begin{aligned} s_3(a) &= \max \{s_2(a) \times p(a|a) \times p(o_1|a), s_2(b) \times p(b|a) \times p(o_1|b), \\ &\quad s_2(c) \times p(c|a) \times p(o_1|c)\} \\ &= \max \{0.0224 \times 0.1 \times 0.2, 0.0192 \times 0.7 \times 0.4, 0.02 \times 0.8 \times 0.2\} \end{aligned}$$

$$s_3(a) = 0.00576$$

$$\begin{aligned} s_3(b) &= \max \{s_2(a) \times p(a|b) \times p(o_1|a), s_2(b) \times p(b|b) \times p(o_1|b), \\ &\quad s_2(c) \times p(c|b) \times p(o_1|c)\} \end{aligned}$$

$$s_3(b) = \max \{0.0224 \times 0.5 \times 0.2, 0.0192 \times 0.2 \times 0.4, 0.02 \times 0.1 \times 0.2\}$$

$$s_3(b) = 0.00448$$

$$\begin{aligned} s_3(c) &= \max \{s_2(a) \times p(a|c) \times p(o_1|a), s_2(b) \times p(b|c) \times p(o_1|b), \\ &\quad s_2(c) \times p(c|c) \times p(o_1|c)\} \\ &= \max \{0.0224 \times 0.8 \times 0.2, 0.0192 \times 0.1 \times 0.4, 0.02 \times 0.8 \times 0.5\} \end{aligned}$$

$$s_3(c) = 0.00768.$$

(4) The most probable final hidden state is the state with the highest viterbi probability at time step  $t = 3$ , in this case, the most probable final hidden state is  $c$ , with a viterbi probability of 0.00768.

9 Suppose we have a chicken from which we collect eggs at noon every day. Now whether or not the chicken has laid eggs for collection depends on some unknown factors that are hidden. We can however assume that the chicken is always in one of 2 states that influence whether the chicken lays eggs, and that this state only depends on the state on the previous day. Now we don't know the state at the initial starting point, we don't know the transition probabilities between the 2 states and we don't know the probability that the chicken lays an egg given a particular state. To start we first guess the transition and emission matrices.

Transition      Emission      Initial

	state 1	state	No eggs	Eggs	state 1	0.2	
state 1	0.5	0.5	State 1	0.3	0.7	State 2	0.8
state 2	0.3	0.7	State 2	0.8	0.2		

We then take a set of observations ( $E = \text{eggs}$ ;  $N = \text{no eggs}$ ):

$N, N, N, N, N, E, E, N, N, N$

This gives us a set of observed transitions between days:

$NN, NN, NN, NN, NE, EC, EN, NN, NN$

using the above information, compute the values in the following tables:

Observed sequence	Highest probability for $S_2 > S_1$	Highest probability of observing that sequence	Prob here	Sequence here
NN	<del><math>S_1 \rightarrow S_2</math></del>	-	0.3804	$S_2, S_2$
NNN	<del><math>S_1 \rightarrow S_2</math></del>	-	0.2408	$S_2, S_2, S_2$
NNNN	<del><math>S_1 \rightarrow S_2</math></del>	-	0.1526	$S_2, S_2, S_2, S_2$
NNNNN	<del><math>S_1 \rightarrow S_2</math></del>	-	0.0968	$S_2, S_2, S_2, S_2, S_2$

NNNNNE	5th N to <sup>1st E</sup> Highest $S_2 \rightarrow S_1$	0.0153	$S_2, S_2, S_2, S_2, S_2, S_1$
NNNNEEE	Highest $S_2 \rightarrow S_1$	0.0048	$S_2, S_2, S_2, S_2, S_2, S_1, S_1$
NN <sup>N</sup> NNEN	Highest $S_2 \rightarrow S_1$	0.0077	$S_2, S_2, S_2, S_2, S_2, S_1, S_1, S_2$
NNNNNEEN	Highest $S_2 \rightarrow S_1$	0.0052	$S_2, S_2, S_2, S_2, S_2, S_1, S_1, S_2$
NNNNNGENNN	$S_1 \rightarrow S_2$	0.00398	$S_2, S_2, S_2, S_2, S_2, S_1, S_1, S_2$
Total	$S_1 \rightarrow S_2$	0	0.00398

New emission matrix estimate:

Observed sequence	Highest probability of observing the sequence if E assumed to come from $S_2$	Highest probability of observing that sequence here
NN	0 (E never comes from $S_2$ )	0.3804
NNN	0 (E never comes from $S_2$ )	0.2408
NNNN	0 (E never comes from $S_2$ )	0.1526
NNNN <del>E</del>	0 (E never comes from $S_2$ )	0.0968
NNN <sup>N</sup> EE <del>E</del>	0 (E never comes from $S_2$ )	0.0153
NNNN <del>N</del> EE <del>E</del>	0 (E never comes from $S_2$ )	0.0048
NNNNNEEN <del>E</del>	0 (E never comes from $S_2$ )	0.0077
NNNNNEENN <del>E</del>	0 (E never comes from $S_2$ )	0.0052
NNNNNGENNN	0 (E never comes from $S_2$ )	0.00398
Total	0	0.00398

10. State the following algorithms in terms of initialization, induction / recursion and termination criteria for HMM in terms of alpha, beta, gamma and di-gamma probabilities. Please define symbols for each and draw a diagram for the induction process of the following algorithms:

(i) Forward Algorithm

(ii) Backward Algorithm

(iii) Viterbi Algorithm

(i) Forward Algorithm

compute likelihood till the  $t-1$  timestamp

update likelihood at time  $t$  based on observation  $o_t$

Let  $\alpha_t(j)$  = joint probability of observing observations  $o_1, o_2, \dots, o_t$  and reaching state  $q_j$  at time  $t$ .

$$\alpha_t(j) = P(o_1, o_2, \dots, o_t, q_t = j | \lambda)$$

$$= \sum_{i=1}^N P(o_1, o_2, \dots, o_{t-1}, q_{t-1} = i, q_t = j | \lambda)$$

Simplifying  $o_1, o_2, \dots, o_{t-1} = O_{1, t-1}$

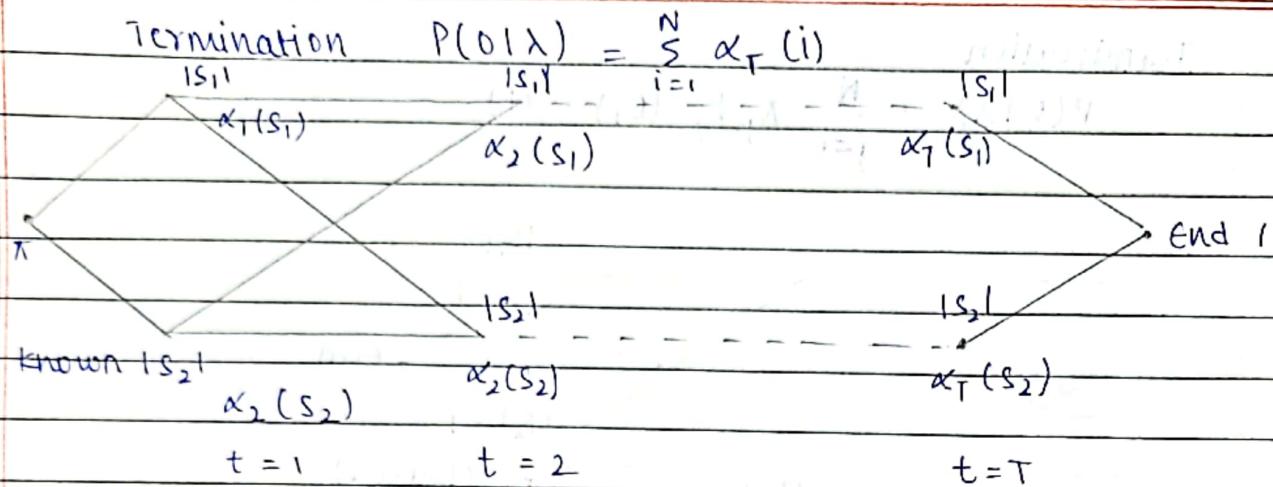
$$= \sum_{i=1}^N P(O_{1, t-1}, q_{t-1} = i) P(q_t = j | q_{t-1} = i, O_{1, t-1}) P(o_t | q_t = j, q_{t-1} = i)$$

$$= \sum_{i=1}^N P(O_{1, t-1}, q_{t-1} = i) \underbrace{P(q_t = j | q_{t-1} = i)}_{\text{independent of } O_{1, t-1}} \underbrace{P(o_t | q_t = j, q_{t-1} = i)}_{\text{independent of } O_{1, t-1}}$$

$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(o_t)$$

$$\text{Recursion } \alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(o_t)$$

$$\text{Initialization } \alpha_1(i) = \pi_i b_i(o_1)$$



## (b) Backward Algorithm

Likelihood of observing  $O$  from  $t+1$  to  $T$  given  $\lambda$

let  $\beta_t(j)$  be probability that future observations  $O_{t+1}, O_{t+2}, \dots, O_T$  have been observed given that hidden state at  $t$  is  $j$

$$\beta_t(j) = P(O_{t+1}, O_{t+2}, \dots, O_T | q_t = j, \lambda)$$

$$\beta_L(j) = \sum_{i=1}^N P(O_{t+1}, O_{t+2}, \dots, O_T, q_{t+1} = i | q_t = j, \lambda)$$

$$\beta_T(j) = \sum_{i=1}^N P(O_{t+1}, O_{t+2}, \dots, O_T, q_{t+1} = i | q_t = j, \lambda)$$

$$\beta_F(j) = \sum_{i=1}^N P(O_{t+1}, O_{t+2}, \dots, O_T | q_{t+1} = i, q_t = j, \lambda) P(q_{t+1} = i | q_t = j, \lambda)$$

$$\beta_T(j) = \sum_{i=1}^N P(O_{t+1}, O_{t+2}, \dots, O_T | q_{t+1} = i, q_t = j, \lambda) a_{ji}$$

$$\beta_t(j) = \sum_{i=1}^N P(O_{t+1}, \dots, O_T | q_{t+1} = i, q_t = j, \lambda) P(O_{t+1} | q_{t+1} = i, q_t = j) a_{ji}$$

$$\beta_t(j) = \sum_{i=1}^N \beta_{t+1}(i) b_i(O_{t+1}) a_{ji}$$

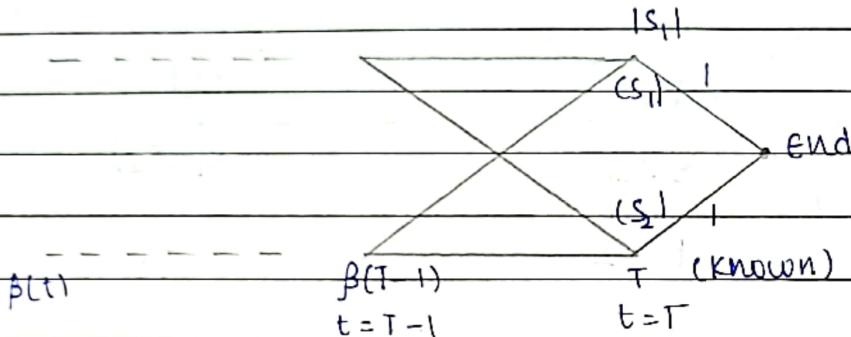
Initialization

$$\beta_1(i) = 1$$

$$\text{Recursion } \beta_T(i) = \sum_{j=1}^N a_{ij} b_j(O_{t+1}) (\beta_{t+1}(j))$$

### Termination

$$P(D | \lambda) = \sum_{j=1}^N \pi_j b_j(O_1) \beta_j(j)$$



### (c) Viterbi Algorithm

$$v_t(j) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_1, q_2, \dots, q_{t-1}, o_1, o_2, \dots, o_{t-1}; q_t=j | O_t | \lambda)$$

### Initialization

$$v_1(j) = \pi_j b_j(O_1)$$

$$b_{t,1}(j) = 0$$

### Recursion

$$v_t(j) = \max_{i=1}^N v_{t-1}(i) \alpha_{ij} b_j(O_t)$$

$$b_{t,1}(j) = \arg \max_{i=1}^N v_{t-1}(i) \alpha_{ij} b_j(O_t)$$

### Termination

$$\text{Best score } P* = \max_{i=1}^N v_T(i)$$

$$\text{Best path start : } q_{T*} = \arg \max_{i=1}^N v_T(i)$$