



50.007 Machine Learning, Fall 2021
Homework 4

Due Tuesday 30 November 2021, 5pm

This homework will be graded by Zhang Qi

In this homework, we would like to look at the Hidden Markov Model (HMM), one of the most influential models used for structured prediction in machine learning.

1. (10 pts) Assume that we have the following training data available for us to estimate the model parameters:

$$\begin{aligned} \text{Count}(X) &= 5 \\ \text{Count}(Y) &= 5 \\ \text{Count}(Z) &= 5 \end{aligned}$$

State sequence	Observation sequence
(X, Y, Z, X)	(b, c, a, b)
(X, Z, Y)	(a, b, c)
(Z, Y, X, Z, Y)	(b, c, a, b, c)
(Z, X, Y)	(c, b, a)

Clearly state what are the parameters associated with the HMM. Under the maximum likelihood estimation (MLE), what would be the values for the optimal model parameters? Clearly show how each parameter is estimated exactly.

2. (10 pts) Now, consider during the evaluation phase, you are given the following new observation sequence. Using the parameters you just estimated from the data, find the most probable state sequence using the Viterbi algorithm discussed in class. Clearly present the steps that lead to your final answer.

State sequence	Observation sequence
(?, ?)	(b, c)

follow that of the lecture slides

3. (10 pts) The Viterbi algorithm discussed in class can be used for computing the single most probable state sequence for a new observation sequence. Specifically, in the Viterbi algorithm we are interested in finding the optimal sequence using the following formula:

$$(s_1^*, s_2^*) = \arg \max_{s_1, s_2} P(s_1, s_2 | o_1 = \mathbf{b}, o_2 = \mathbf{c})$$

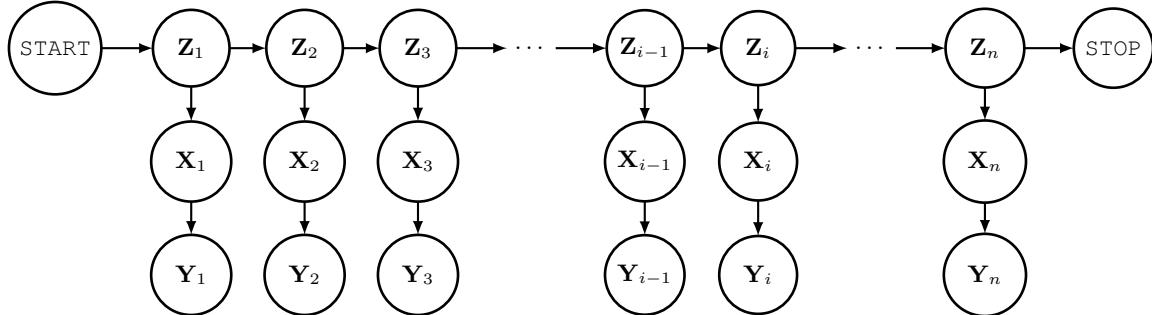
However, an alternative method to estimate the state sequence is finding the optimal state of index i by solving the optimization problem

joints probability

$$s_i^* = \arg \max_{s_i} P(s_i | o_1 = \mathbf{b}, o_2 = \mathbf{c})$$

where $P(s_i | o_1 = \mathbf{b}, o_2 = \mathbf{c})$ is the marginal distribution and $i \in \{1, 2\}$. Please use the forward and backward scores for HMM discussed in class to calculate the marginal distribution $P(s_i | o_1 = \mathbf{b}, o_2 = \mathbf{c})$. Clearly present the steps that lead to your final answer.

4. (20 pts) Now consider a slightly different graphical model which extends the HMM (see below). For each state (Z), there is now an observation pair (X, Y), where Y sequence is generated from the X sequence.



Assume you are given a large collection of observation pair sequence, and a predefined set of possible states, you would like to estimate the most probable state sequence for each observation pair sequence using an EM algorithm similar to the dynamic programming algorithm discussed in class. Clearly define the forward and backward scores in a way analogous to those defined for HMM, and explain what they mean. Give algorithms for computing the forward and backward scores. Analyze the complexity associated with your algorithms.

1. Two parameters:

① Emission probability $[a_{u,v}]$

↳ output obtain from state (probability)

② Transition probability $[b_{u(O)}]$

↳ The probability of moving from one state to the next state.

$$a_{u,v} = \frac{\text{Count}(u,v)}{\text{Count}(u)}$$

$$b_{u(O)} = \frac{\text{Count}(u \rightarrow O)}{\text{Count}(u)}$$

• Transition probability

Sanity check (Σ)

u v	x	y	z	STOP
START	$\frac{1}{2}$	0	$\frac{1}{2}$	0
x	0	$\frac{2}{5}$	$\frac{2}{5}$	$\frac{1}{5}$
y	$\frac{1}{5}$	0	$\frac{1}{5}$	$\frac{3}{5}$
z	$\frac{2}{5}$	$\frac{3}{5}$	0	0

, , , ,

The table passes the sanity check.

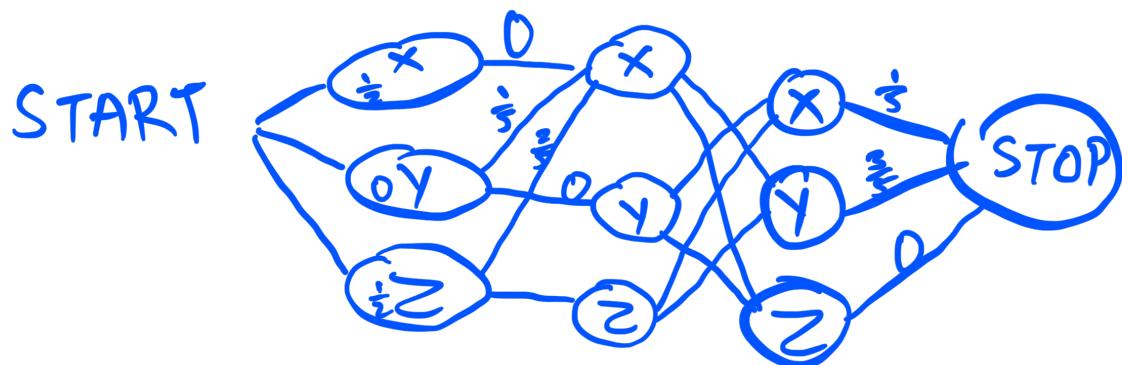
• Emission probability

Sanity check (Σ)

u o	a	b	c
x	$\frac{2}{5}$	$\frac{3}{5}$	$\frac{0}{5}$
y	$\frac{1}{5}$	$\frac{0}{5}$	$\frac{4}{5}$
z	$\frac{1}{5}$	$\frac{3}{5}$	$\frac{1}{5}$

, , , ,

Q2. Using the backtracking algorithm,



- ① We know that we cannot move from START to state Y ($a_{\text{START}, Y} = 0$)
- ② Our state before STOP cannot be (Z)

Via Backtracking, we trace a few paths:

ⓐ $\text{START} \rightarrow X \rightarrow Y \rightarrow \text{STOP}$

$$\begin{aligned}\text{likelihood} &= a_{\text{START}, X} + a_{X, Y} + a_{Y, \text{STOP}} \times \\ &\quad b_X(b) \times b_Y(c) \\ &= \frac{1}{2} \times \frac{2}{5} \times \frac{3}{5} \times \frac{3}{5} \times \frac{4}{5} \\ &= 0.0576\end{aligned}$$

ⓑ $\text{START} \rightarrow Z \rightarrow Y \rightarrow \text{STOP}$

$$\begin{aligned}\text{likelihood} &= a_{\text{START}, Z} + a_{Z, Y} + a_{Y, \text{STOP}} \times \\ &\quad b_Z(b) \times b_Y(c) \\ &= \frac{1}{2} \times \frac{3}{5} \times \frac{3}{5} \times \frac{3}{5} \times \frac{4}{5} \\ &= 0.0864\end{aligned}$$

Path ⓑ seems short (But not a very gd way to find)

Initialisation $\pi(0, n) = \begin{cases} 1 & \text{if } n = \text{START} \\ 0 & \text{otherwise} \end{cases}$

$$\textcircled{1} \quad \pi(1, X) = a_{\text{START}, X} \times b_X(b) = \frac{1}{2} \times \frac{3}{5} = \frac{3}{10}$$

$$\pi(1, Z) = a_{\text{START}, Z} \times b_Z(b) = \frac{1}{2} \times \frac{3}{5} = \frac{3}{10}$$

$$\pi(1, Y) = a_{\text{START}, Y} \times b_Y(b) = 0$$

$$\textcircled{2} \quad \pi(2, X) = \max \{ \pi(1, u) \times a_{u, X} \times b_X(c) \}$$

$$= \max \{ \pi(1, X) \times a_{X, X} \times b_X(c),$$

$$\pi(1, Z) \times a_{Z, X} \times b_X(c),$$

$$\pi(1, Y) \times a_{Y, X} \times b_X(c) \}$$

$$= \max \{ 0, 0, 0 \}$$

$$= 0$$

$$\pi(2, Y) = \max \{ \pi(1, u) \times a_{u, Y} \times b_Y(c) \}$$

$$= \max \{ \pi(1, X) \times a_{X, Y} \times b_Y(c),$$

$$\pi(1, Y) \times a_{Y, Y} \times b_Y(c),$$

$$\pi(1, Z) \times a_{Z, Y} \times b_Z(c) \}$$

$$= \max \{ \frac{3}{10} \times \frac{2}{5} \times \frac{4}{5}, 0, \frac{3}{10} \times \frac{3}{5} \times \frac{4}{5} \}$$

$$= \max \{ 0.096, 0, 0.144 \}$$

$$= 0.144 \#$$

$$\begin{aligned}
 \pi(2, z) &= \max \{ \pi(1, u) \times a_{u,z} \times b_z(c) \} \\
 &= \max \{ \pi(1, x) \times a_{x,z} \times b_z(c) \\
 &\quad \pi(1, y) \times a_{y,z} \times b_z(c) \\
 &\quad \pi(1, z) \times a_{z,z} \times b_z(c) \\
 &\quad \} \\
 &= \max \{ \frac{3}{10} \times \frac{2}{5} \times \frac{1}{3}, 0, 0 \} \\
 &= 0.024 \text{ #}
 \end{aligned}$$

$$\begin{aligned}
 ③ \pi(3, \text{STOP}) &= \max \{ \pi(2, \pi) \times a_{\pi, \text{STOP}} \} \\
 &= \max \{ \pi(2, \pi) \times a_{x, \text{STOP}}, \\
 &\quad \pi(2, \pi) \times a_{y, \text{STOP}}, \\
 &\quad \pi(2, \pi) \times a_{z, \text{STOP}} \\
 &\quad \} \\
 &= \max \{ 0, 0.144 \times \frac{3}{5}, 0 \} \\
 &= 0.0864 \text{ #}
 \end{aligned}$$

∴ Final Path

START → Z → Y → STOP

3. From leftmost notes 4,

$$s_i^* = \arg \max_{s_i} P(s_i | o_1 = b, o_2 = c)$$

↓ relate to the problem

To find the optimal tag at the position i given the entire observation sequence of (o_1, o_2)

We can then apply the result,

$$s_i^* = \arg \max_u (X_u(i) \beta_u(i)), i \in \{1, 2\}$$

Using forward probability, we can define $X_u(i)!$

/ Base case:

$$X_u(1) = X_{\text{START}, u}$$

Recursive:

$$X_u(i+1) = \sum_j X_u(j) a_{j,u} b_u(x_j) \quad \begin{array}{l} \text{Stop @} \\ j+1=3 \end{array}$$

$$\begin{aligned} X_x(1) &= X_{\text{START}}, x = \frac{1}{2} \\ X_y(1) &= X_{\text{START}}, y = 0 \\ X_z(1) &= X_{\text{START}}, z = \frac{1}{2} \end{aligned} \quad \begin{array}{l} \text{Base} \\ \text{case:} \end{array}$$

$$\begin{aligned} X_x(2) &= X_x(1) \times a_{x,x} \times b_x(b) + \\ &\quad X_y(1) \times a_{y,x} \times b_y(b) + \\ &\quad X_z(1) \times a_{z,x} \times b_z(b) \\ &= 0 + 0 + \left(\frac{1}{2} \times \frac{2}{3} \times \frac{3}{2}\right) \\ &= \frac{3}{25} \\ &= 0.12 \end{aligned}$$

$$\begin{aligned} X_y(2) &= X_x(1) \times a_{x,y} \times b_x(b) + \\ &\quad X_y(1) \times a_{y,y} \times b_y(b) + \\ &\quad X_z(1) \times a_{z,y} \times b_z(b) \\ &= \left(\frac{1}{2} \times \frac{2}{3} \times \frac{2}{3}\right) + 0 + \left(\frac{1}{2} \times \frac{1}{3} \times \frac{2}{3}\right) \\ &= 0.3 \end{aligned}$$

$$\begin{aligned}
 K_x(2) &= K_x(1) \times a_{x,z} \times b_x(b) + \\
 K_y(1) \times a_{y,z} \times b_y(b) + \\
 K_z(1) \times a_{z,z} \times b_z(b) + \\
 &= \frac{1}{2} \times 0.4 \times 0.6 + 0 + 0 \\
 &= 0.12
 \end{aligned}$$

For Backward Score calculation,

Base case:

$$\beta_u(n) = a_{u,\text{STOP}} b_u(x_n)$$

Recursive case:

$$\beta_u(j) = \sum_v a_{u,v} b_u(x_j) \beta_v(j+1)$$

$$\beta_x(2) = a_{x,\text{STOP}} b_x(c) = 0$$

$$\beta_y(2) = a_{y,\text{STOP}} b_y(c) = 0.6 \times \frac{4}{5} = 0.48$$

$$\beta_z(2) = a_{z,\text{STOP}} b_z(c) = 0$$

$$\begin{aligned}
 \beta_x(1) &= a_{x,x} b_x(b) \times \beta_x(2) + \\
 &\quad a_{x,y} b_x(b) \times \beta_y(2) + \\
 &\quad a_{x,z} b_x(b) \times \beta_z(2) \\
 &= 0 + \frac{2}{5} \cdot \frac{3}{5} \cdot 0.48 + 0 \\
 &= 0.1152
 \end{aligned}$$

$$\begin{aligned}
 \beta_y(2) &= a_{y,x} b_y(b) \times \beta_x(2) + \\
 &\quad a_{y,y} b_y(b) \times \beta_y(2) + \\
 &\quad a_{y,z} b_y(b) \times \beta_z(2) \\
 &= 0 + 0 + 0 \\
 &= 0
 \end{aligned}$$

$$\begin{aligned}
 \beta_z(2) &= a_{z,x} b_z(b) \times \beta_x(2) \\
 &\quad a_{z,y} b_z(b) \times \beta_y(2) \\
 &\quad a_{z,z} b_z(b) \times \beta_z(2) \\
 &= 0 + \frac{3}{5} \times \frac{3}{5} \times 0.48 + 0 \\
 &= 0.1728
 \end{aligned}$$

$$s_1^* = \arg \max_u \{ \alpha_u(1) \cdot \beta_u(1) \}$$

$$= \arg \max_u \{ K_x(1) \cdot \beta_x(1), K_y(1) \cdot \beta_y(1), K_z(1) \cdot \beta_z(1) \}$$

$$= \arg \max_u \{ 0.0576, 0, 0.0804 \}$$

$$= 2$$

$$s_2^* = \arg \max_u (K_u(2) \cdot \beta_u(2))$$

$$= \arg \max_u \{ K_x(2) \times \beta_x(2), K_y(2) \times \beta_y(2), K_z(2) \times \beta_z(2) \}$$

$$= \arg \max_u \{ 0, 0.144, 0 \}$$

$$= 1$$

Q 4.
Using soft EM, we can use the results to derive the following for our use.

$$a_{v,u} = \frac{\text{count}(v,u)}{\text{count}(v)} \quad b_v(x_i) = \frac{\text{count}(v \rightarrow x_i)}{\text{count}(v)}$$

$$c_v(y_i) = \frac{\text{count}(v \rightarrow y_i)}{\text{count}(v)}$$

\downarrow
how
emission
for each
layer.

∴ Using the above, we can compute the forward and backward scores respectively.

- Forward score,
 $K_u(i)$ is the sum of scores of all paths from start to state (i,u)

$$a_{u,i} = a_{\text{start},u}$$

$$a_{u,i+1} = \sum_v K_v(i) a_{v,u} b_v(x_i) c_v(y_i) \text{ for } u \in 1, \dots, N-1, i = 1, \dots, n-1$$

- Backward score,

$\beta_n(i)$ is the sum of scores of all paths from state (i, u) to STOP

Base case:

$$\beta_n(n) = a_{n, \text{STOP}} b_n(x_n) \cdot c_n(y_n) \quad \forall n \in 1, \dots, N-1$$

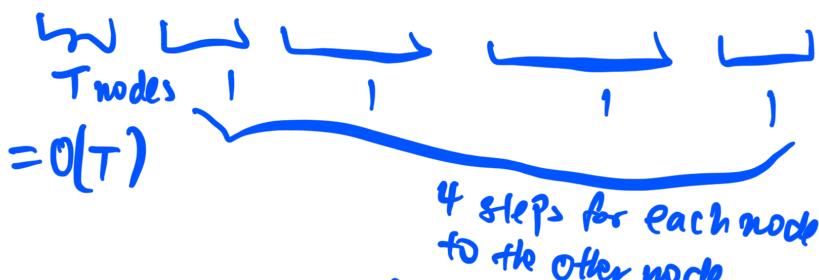
Recursive Case:

$$\beta_n(i) = \sum_v a_{i,v} b_n(x_i) \beta_v(i+1) \cdot c_n(y_i) \quad \forall n \in 1, \dots, N-1, i = n-1, \dots, 1$$

For each layer, $O(T)$ time complexity!

$\beta_n(i)$ is $O(T)$ time complexity.

$$\beta_n(i) = \sum_v (a_{i,v} \cdot b_n(x_i) \cdot c_n(y_i) \cdot \beta_v(i+1))$$



$$O(T^2 n) = O(T) \cdot O(T) \cdot O(n)$$

Explanation: ① for each node to connect to T other nodes in computation ($O(T)$)

② for all T nodes in each layer to connect to T other nodes, $(O(T) \times O(T))$ in computation

For n layers, to repeat ① and ② in each layer,

$$\begin{aligned} \text{Total time complexity} &= O(T) \times O(T) \times O(n) \\ &= O(T^2 n) \end{aligned}$$

* We discount 4 as constants do not matter too much in time complexity
(Refer to introduction to Algorithms (50-004))

For space complexity, we store T values in each layer.

∴ For n layers,

$$\text{Space complexity} = O(T \times n) = O(Tn)$$