

50.007 Machine Learning, Fall 2015

Homework 4

Due Monday 23 Nov 2015, 5pm

Sample Solutions

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.

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

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

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.

Answer The transition probabilities are estimated as:

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

	X	Y	Z	STOP
START	1/2	0	1/2	0
X	2/7	0	3/7	2/7
Y	1/4	0	0	3/4
Z	1/7	4/7	1/7	1/7

The emission probabilities are estimated as:

$$b_u(o) = \frac{\text{Count}(u \rightarrow o)}{\text{Count}(u)}$$

	a	b	c	d
X	3/7	2/7	2/7	0
Y	1/2	0	1/4	1/4
Z	1/7	4/7	1/7	1/7

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
$(?, ?)$	(\mathbf{a}, \mathbf{d})

Answer

- Base case:

$$\pi(0, \text{START}) = 1, \quad \text{otherwise } \pi(0, v) = 0 \text{ if } v \neq \text{START} \quad (1)$$

- Moving forward:

$$k = 1$$

$$\pi(1, X) = a_{\text{START}, X} \times b_X(a) = 1 \times 1/2 \times 3/7 = 3/14 \quad (2)$$

$$\pi(1, Y) = a_{\text{START}, Y} \times b_Y(a) = 0 \times 1/2 = 0 \quad (3)$$

$$\pi(1, Z) = a_{\text{START}, Z} \times b_Z(a) = 1 \times 1/2 \times 1/7 = 1/14 \quad (4)$$

$$k = 2$$

$$\begin{aligned} \pi(2, X) &= \max_{u \in \mathcal{T}} \{ \pi(1, u) \times a_{u, X} \times b_X(d) \} \\ &= \max \{ 3/14 \times 2/14 \times 0, \quad 0 \times 1/4 \times 0, \quad 1/4 \times 1/7 \times 0 \} \\ &= 0 \end{aligned} \quad (5)$$

$$\begin{aligned} \pi(2, Y) &= \max_{u \in \mathcal{T}} \{ \pi(1, u) \times a_{u, Y} \times b_Y(d) \} \\ &= \max \{ 3/14 \times 0 \times 1/4, \quad 0 \times 0 \times 1/4, \quad 1/14 \times 4/7 \times 1/4 \} \\ &= 1/98 \end{aligned} \quad (6)$$

$$\begin{aligned} \pi(2, Z) &= \max_{v \in \mathcal{T}} \{ \pi(1, v) \times a_{v, Z} \times b_Z(d) \} \\ &= \max \{ 3/14 \times 3/7 \times 1/7, \quad 0 \times 0 \times 1/7, \quad 1/14 \times 1/7 \times 1/7 \} \\ &= 9/686 \end{aligned} \quad (7)$$

$$k = 3$$

$$\begin{aligned} \pi(3, \text{STOP}) &= \max_{v \in \mathcal{T}} \{ \pi(2, v) \times a_{v, \text{STOP}} \} \\ &= \max \{ 0 \times 2/7, 1/98 \times 3/4, 9/686 \times 1/7 \} \\ &= 3/392 \end{aligned} \quad (8)$$

- Backtracking:

$$y_2^* = \arg \max_{v \in \mathcal{T}} \{\pi(2, v) \times a_{v, \text{STOP}}\} = Y \quad (9)$$

$$y_1^* = \arg \max_{v \in \mathcal{T}} \{\pi(1, v) \times a_{v, Y}\} = Z \quad (10)$$

Therefore, the optimal sequence is: Z, Y .

3. (20 pts) The HMM discussed in class makes a simple first-order assumption, where the next state only depends on the previous state in the generative process. However, it is possible to extend the model discussed in class to have second-order dependencies. In other words, the HMM can be parameterised in the following way:

$$p(x_1, \dots, x_n, y_1, y_2, \dots, y_n) = \prod_{i=1}^{n+1} p(y_i | y_{i-2}, y_{i-1}) \times \prod_{i=1}^n p(x_i | y_i)$$

where we define $y_{-1} = y_0 = \text{START}$ and $y_{n+1} = \text{STOP}$.

In other words, the transition probabilities are changed from $p(y_i | y_{i-1})$ to $p(y_i | y_{i-2}, y_{i-1})$ now. Describe the Viterbi algorithm used for decoding such a second-order HMM model. In other words, describe the dynamic programming algorithm that computes the following efficiently for such an HMM:

$$(y_1^*, y_2^*, \dots, y_n^*) = \arg \max_{y_1, y_2, \dots, y_n} p(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n)$$

Answer The key idea is to consider a pair of states $\langle u, v \rangle$ instead of a single state v at each time step, where v is the current state, and u is the last state (the state at the previous time step).

Define $a_{t,u,v} \equiv p(y_i = v | y_{i-2} = t, y_{i-1} = u)$, and $b_u(o) \equiv p(x_i = o | y_i = u)$.

- Base cases:

$$\pi(0, \langle \text{START}, \text{START} \rangle) = 1 \quad (11)$$

$$\pi(0, \langle u, v \rangle) = 0, \text{ if } \langle u, v \rangle \neq \langle \text{START}, \text{START} \rangle \quad (12)$$

$$\pi(1, \langle \text{START}, v \rangle) = a_{\text{START}, \text{START}, v} \times b_v(x_1) \quad \forall v \quad (13)$$

$$\pi(1, \langle u, v \rangle) = 0, \text{ if } u \neq \text{START} \quad (14)$$

- Moving forward recursively

For $k \geq 2$

$$\pi(k, \langle u, v \rangle) = \max_{t \in \mathcal{T}} \{\pi(k-1, \langle t, u \rangle) \times a_{t,u,v} \times b_v(x_k)\} \quad (15)$$

- Final transition

From $\langle y_{n-1}, y_n \rangle$ to STOP, the transition:

$$\pi(n+1, \langle u, \text{STOP} \rangle) = \max_{t \in \mathcal{T}} \{\pi(n, \langle t, u \rangle) \times a_{t,u, \text{STOP}}\} \quad (16)$$

- Backward tracking

$$y_n^* = \arg \max_u \{ \pi(n+1, \langle u, \text{STOP} \rangle) \} \quad (17)$$

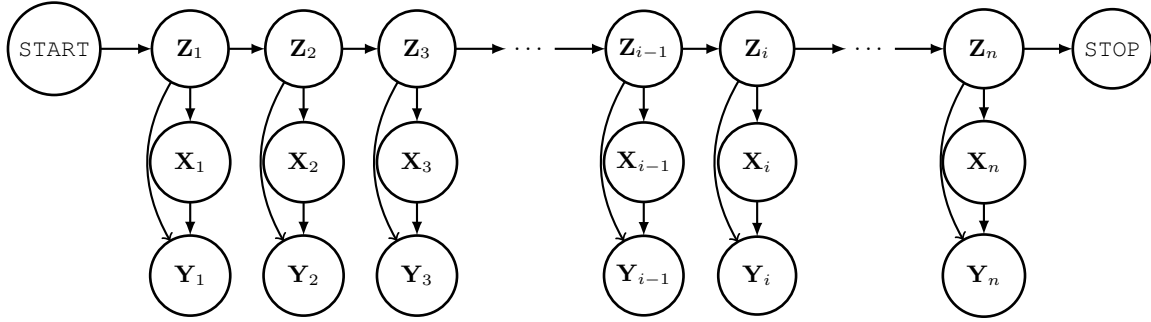
$$\begin{aligned} y_k^* &= \arg \max_t \{ \pi(k+1, \langle t, y_{k+1}^* \rangle) \times a_{t, y_{k+1}^*, y_{k+2}^*} \times b_{y_{k+2}^*}(x_{k+2}) \} \\ &= \arg \max_t \{ \pi(k+1, \langle t, y_{k+1}^* \rangle) \times a_{t, y_{k+1}^*, y_{k+2}^*} \} \end{aligned} \quad (18)$$

Note: alternatively, if you do not compute $\pi(n+1, \langle u, \text{STOP} \rangle)$, that is also fine. You can obtain the last two optimal states from the following:

$$\langle y_{n-1}^*, y_n^* \rangle = \arg \max_{u,v} \{ \pi(n, \langle u, v \rangle) \times a_{u,v,\text{STOP}} \} \quad (19)$$

Why is it so? Think about it! Also think about whether this alternative approach would be faster!

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



Assume you are given a large collection of observation pair sequences, and a predefined set of possible states $\{0, 1, \dots, N-1, N\}$, where $0 = \text{START}$ and $N = \text{STOP}$. You would like to estimate the most probable state sequence for each observation pair sequence using an algorithm similar to the dynamic programming algorithm discussed in class. Clearly define the forward and backward scores in a way analogous to HMM. Give algorithms for computing the forward and backward scores. Analyze the time complexity associated with your algorithms (for an observation pair sequence of length n).

Answer Assume we have a set of possible states $\{0, 1, \dots, N-1, N\}$ where $0 = \text{START}$ and $N = \text{STOP}$.

$$\begin{aligned} &P(x_1, \dots, x_{i-1}, y_1, \dots, y_{i-1}, z_i = u, x_i, \dots, x_n, y_i, \dots, y_n; \theta) \\ &= P(x_1, \dots, x_{i-1}, y_1, \dots, y_{i-1}, z_i = u; \theta) \times P(x_i, \dots, x_n, y_i, \dots, y_n | z_i = u; \theta) \\ &= \alpha_u(i) \beta_u(i) \end{aligned} \quad (20)$$

where

$$\alpha_u(i) = P(x_1, \dots, x_{i-1}, y_1, \dots, y_{i-1}, z_i = u; \theta) \quad (21)$$

$$\beta_u(i) = P(x_i, \dots, x_n, y_i, \dots, y_n | z_i = u; \theta) \quad (22)$$

Forward

- Base Case

$$\alpha_u(1) = a_{\text{START},u}, \quad \forall u \in \{1, \dots, N-1\} \quad (23)$$

- Moving forward

For $i = 2, \dots, n$:

$$\alpha_u(i+1) = \sum_v \alpha_v(i) \times a_{v,u} \times b_v(x_i) \times c_{v,x_i}(y_i) \quad (24)$$

where

$$c_{v,x}(y) = P(y|v, x) \quad (25)$$

Backward

- Base case

$$\beta_u(n) = a_{u,\text{STOP}} \times b_u(x_n) \times c_{u,x_n}(y_n) \quad \forall u = 1, \dots, N-1 \quad (26)$$

- Moving forward

For $i = n-1, \dots, 1$:

$$\beta_u(i) = \sum_v a_{u,v} \times b_u(x_i) \times c_{u,x_i}(y_i) \times \beta_v(i+1) \quad (27)$$

At each time step/position, there are N forward (α) and N backward (β) terms to compute. To compute each term, there are $O(N)$ operations. Thus, at each time step/position, there are $O(N^2)$ operations. The length of sentence is n , which is the number of different time steps/positions. Hence, the total complexity is $O(nN^2)$.