

01.112 Machine Learning

Homework 4

Joel Huang, 1002530

November 23, 2018

Hidden Markov Models

Model Parameters

Parameters associated with the HMM

The HMM is parameterized by:

1. \mathcal{T} , the set of states, including the **START** and **STOP** states. $\mathcal{T} = 0$ (**START**), \dots , $|\mathcal{T}| - 2$ (**STOP**).
2. \mathcal{O} , the set of observation symbols.
3. $a_{u,v}$, the transition parameters, which are the probabilities of transitioning from state u to v .

$$a_{u,v} = p(y_{next} = v \mid y_{curr} = u) \quad (1)$$

where $u \in [0, \dots, |\mathcal{T}| - 2]$, $v \in [1, \dots, |\mathcal{T}| - 1]$.

4. $b_u(o)$, the emission parameters, which are the probabilities of emitting symbol o given state u .

$$b_u(o) = p(x = o \mid y = u) \quad (2)$$

Computing optimal model parameters

- States $\mathcal{T} = \{\text{START}, X, Y, Z, \text{STOP}\}$
- Observations $\mathcal{O} = \{a, b, c\}$
- Transition parameters:

$$a_{\text{START},X} = \frac{\text{Count}(\text{START}; X)}{\text{Count}(\text{START})} = \frac{2}{4} = 0.5 \quad (3)$$

$$a_{\text{START},Z} = \frac{\text{Count}(\text{START}; Z)}{\text{Count}(\text{START})} = \frac{2}{4} = 0.5 \quad (4)$$

$$a_{X,Y} = \frac{\text{Count}(X; Y)}{\text{Count}(X)} = \frac{2}{5} = 0.4 \quad (5)$$

$$a_{X,Z} = \frac{\text{Count}(X; Z)}{\text{Count}(X)} = \frac{2}{5} = 0.4 \quad (6)$$

$$a_{Y,X} = \frac{\text{Count}(Y; X)}{\text{Count}(Y)} = \frac{1}{5} = 0.2 \quad (7)$$

$$a_{Y,Z} = \frac{\text{Count}(Y; Z)}{\text{Count}(Y)} = \frac{1}{5} = 0.2 \quad (8)$$

$$a_{Z,X} = \frac{\text{Count}(Z; X)}{\text{Count}(Z)} = \frac{2}{5} = 0.4 \quad (9)$$

$$a_{Z,Y} = \frac{\text{Count}(Z; Y)}{\text{Count}(Z)} = \frac{3}{5} = 0.6 \quad (10)$$

$$a_{X,\text{STOP}} = \frac{\text{Count}(X; \text{STOP})}{\text{Count}(X)} = \frac{1}{5} = 0.2 \quad (11)$$

$$a_{Y,\text{STOP}} = \frac{\text{Count}(Y; \text{STOP})}{\text{Count}(Y)} = \frac{3}{5} = 0.6 \quad (12)$$

- Emission parameters:

$$b_X(a) = \frac{\text{Count}(X \rightarrow a)}{\text{Count}(X)} = \frac{2}{5} = 0.4 \quad (13)$$

$$b_Y(a) = \frac{\text{Count}(Y \rightarrow a)}{\text{Count}(Y)} = \frac{2}{5} = 0.4 \quad (14)$$

$$b_Z(a) = \frac{\text{Count}(Z \rightarrow a)}{\text{Count}(Z)} = \frac{1}{5} = 0.2 \quad (15)$$

$$b_X(b) = \frac{\text{Count}(X \rightarrow b)}{\text{Count}(X)} = \frac{3}{5} = 0.6 \quad (16)$$

$$b_Z(b) = \frac{\text{Count}(Z \rightarrow b)}{\text{Count}(Z)} = \frac{3}{5} = 0.6 \quad (17)$$

$$b_Y(c) = \frac{\text{Count}(Y \rightarrow c)}{\text{Count}(Y)} = \frac{3}{5} = 0.6 \quad (18)$$

$$b_Z(c) = \frac{\text{Count}(Z \rightarrow c)}{\text{Count}(Z)} = \frac{1}{5} = 0.2 \quad (19)$$

Viterbi Algorithm

Formulation

1. Our objective is to output the sequence of tags y_1, \dots, y_n with the highest likelihood:

$$\max_{y_1, \dots, y_n} P(y_1, \dots, y_n) \quad (20)$$

Let y_1, \dots, y_k be a subset of y_1, \dots, y_n , where $v = y_k$. Define $\pi(k, v)$ as the sequence of tags up to y_k that has the highest likelihood:

$$\pi(k, v) = \max_{y_1, \dots, y_k} P(y_1, \dots, y_k) \quad (21)$$

$$\max_{y_1, \dots, y_k} \left\{ \prod_{i=1}^k a_{y_{i-1}, y_i} \cdot \prod_{i=1}^k b_{y_i}(x_i) \right\}$$

Table of $\pi(k, v)$ for next state v				
iteration k	0	1	...	n
$v = \text{START}$	1	0	...	0
$v = y_1$	0	$\pi(1, y_1)$...	$\pi(n, y_1)$
$v = y_2$	0	$\pi(1, y_2)$...	$\pi(n, y_2)$
...	\vdots	\vdots	\ddots	\vdots
$v = y_n$	0	$\pi(1, y_n)$...	$\pi(n, y_n)$

For example, for iteration $k = 4$, and next state $v = y_2$, we find the maximum probability of the over all possible states.

$$\pi(4, y_2) = \max_{u \in y_0, \dots, y_{n+1}} \{ \pi(3, u) \cdot a_{u, y_2} \cdot b_{y_2}(x_4) \}, \quad (22)$$

- $\pi(3, u)$, where u is the state which survives after taking the maximum probability over all possible states when $k = 3$,
- a_{u, y_2} , the probability of transitioning from state u to state y_2 ,
- $b_{y_2}(x_4)$, the probability of emitting observation x_4 given state y_2 .

2. Base case, $k = 0$.

$$\pi(k = 0, v) = \begin{cases} 1, & \text{if } v = \text{START} \\ 0, & \text{otherwise} \end{cases} \quad (23)$$

3. Recursive, for $k = 1, \dots, n$, the best probability of the best previous path \times transition \times emission.

$$\pi(k, v) = \max_u \{ \pi(k-1, u) \cdot a_{u, v} \cdot b_v(x_k) \} \quad (24)$$

4. Termination case.

$$\pi(k = (n+1), \text{STOP}) = \max_v \{ \pi(n, v) \cdot a_{v, \text{STOP}} \} \quad (25)$$

Computation

1. Input: $x = \{b, b\}$

2. $k = 0$, base case:

$$\pi(0, \text{START}) = 1 \quad (26)$$

3. $k = 1$:

$$\begin{aligned} \pi(1, X) &= \pi(0, \text{START}) \cdot a_{\text{START}, X} \cdot b_X(b) \\ &= 1 \times 0.5 \times 0.6 = 0.3 \\ \pi(1, Y) &= \pi(0, \text{START}) \cdot a_{\text{START}, Y} \cdot b_Y(b) \\ &= 1 \times 0 \times 0 = 0 \\ \pi(1, Z) &= \pi(0, \text{START}) \cdot a_{\text{START}, Z} \cdot b_Z(b) \\ &= 1 \times 0.5 \times 0.6 = 0.3 \end{aligned} \quad (27)$$

4. $k = 2$:

$$\begin{aligned} \pi(2, X) &= \max_{X, Y, Z} \{ \pi(1, X) \cdot a_{X, X} \cdot b_X(b), \\ &\quad \pi(1, Y) \cdot a_{X, Y} \cdot b_Y(b), \\ &\quad \pi(1, Z) \cdot a_{X, Z} \cdot b_Z(b) \} \\ &= \max \{0, 0, 0.072\} = 0.072 \\ \pi(2, Y) &= \max \{0, 0, 0\} = 0 \\ \pi(2, Z) &= \max \{0.072, 0, 0\} = 0.072 \end{aligned} \quad (28)$$

5. $k = 3$, termination case:

$$\begin{aligned} \pi(3, \text{STOP}) &= \max_v \{ \pi(2, v) \cdot a_{v, \text{STOP}} \} \\ &= \max \{ \pi(2, X) \cdot a_{X, \text{STOP}} \\ &\quad \pi(2, Y) \cdot a_{Y, \text{STOP}}, \\ &\quad \pi(2, Z) \cdot a_{Z, \text{STOP}} \} \\ &= \max \{0.072 \times 0.2, 0, 0\} = 0.0144 \end{aligned} \quad (29)$$

6. We then backtrack, finding the best values for each state:

$$\begin{aligned} y_2 &: \arg \max_v \{0.0144, 0, 0\} = X \\ y_1 &: \arg \max_v \{0, 0, 0.072\} = Z \end{aligned} \quad (30)$$

So the most probable sequence is:

$$y_0, \dots, y_{n+1} = \text{START}, Z, X, \text{STOP} \quad (31)$$

Top- k decoding

Currently, at each $\pi(i, u)$ we only store one parent. However suppose we stored the top k predecessors. So each $\pi(i, u)$ corresponds not only to a most likely value and the node which it transitioned from, but a list of the top k nodes it could have transitioned from and their values in sorted order. Therefore, in order to perform top- k decoding, we must store the top k optimal sub-paths at each node, instead of just the top sub-path with the highest probability.

Formulation

1. Base case, $i = 0$. This is unchanged from the vanilla Viterbi algorithm.

$$\pi(i = 0, v) = \begin{cases} 1, & \text{if } v = \text{START} \\ 0, & \text{otherwise} \end{cases} \quad (32)$$

2. Recursive, for $i = 1, \dots, n$, the k best probabilities of the previous paths \times transition \times emission. Here, in contrast to the vanilla Viterbi algorithm which carries out a max operation to sieve out the best preceding node, we take the top k preceding nodes and store them.

We define a k-max operator, which selects the k highest elements in the set, then sorts them in descending order.

$$\pi(i, v) = \text{k-max}_u \{ \pi(i-1, u) \cdot a_{u,v} \cdot b_v(x_i) \} \quad (33)$$

3. Termination case.

$$\pi(i = (n+1), \text{STOP}) = \text{k-max}_v \{ \pi(n, v) \cdot a_{v,\text{STOP}} \} \quad (34)$$

4. Backtracking. We can visualize the backtracking using a matrix A , of size $k \times n$.

$$\begin{array}{c} \hline 1 \quad \dots \quad n \\ 1 \\ \vdots \\ k \\ \hline \end{array}$$

At the termination point, we have $\pi(n, v) \cdot a_{v,\text{STOP}}$ for all $v \in \mathcal{T}$. Since $\pi(n+1, \text{STOP})$ is now a list of the k highest probabilities of paths, we check which paths from the n -th layer led to it. For each v , if $\pi(n, v) \cdot a_{v,\text{STOP}} \in \pi(n+1, \text{STOP})$, then fill the n -th column of A with the states v that contributed to these highest probabilities.

Similarly, for the i -th layer, compute $\pi(i-1, u) \cdot a_{u,v} \cdot b_v(x_i)$ and check which path from the $(i-1)$ -th layer led to it, and fill the i -th column of A with the states u that contributed to the highest probabilities. Repeat this till the first ($i = 1$) layer, where we instead compute $\pi(1, u) \cdot a_{\text{START},u} \cdot b_u(x_1)$. Reading the rows of A then gives us each top- k most probable sequences.

Forward-backward algorithm

Define two observation spaces $\{X_1, \dots, X_n\}$, and $\{Y_1, \dots, Y_n\}$. Define a state space $\{Z_0, \dots, Z_{n+1}\}$ where $Z_0 = \text{START}$, $Z_{n+1} = \text{STOP}$. Parameterize the model by $a_{u,v}$, $b_u(o)$, and $c_o(e)$, where $c_o(e)$ is the additional emission probability introduced. Unfolding the joint probability using newly defined α and β :

$$\begin{aligned} & P(x_1, \dots, x_n, y_1, \dots, y_n, z_i = u) \\ &= P(x_1, \dots, x_{i-1}, y_1, \dots, y_{i-1}, z_i = u) \cdot \\ & \quad P(x_i, \dots, x_n, y_i, \dots, y_n | z_i = u) \\ &= P(x_1, \dots, x_{i-1}, z_i = u) \cdot \sum_{j=1}^i P(y_j | x_j) \cdot \\ & \quad P(x_i, \dots, x_{i-1} | z_i = u) \cdot \sum_{k=i}^n P(y_k | x_k) \\ &= \alpha_u(i) \cdot \beta_u(i) \end{aligned} \quad (35)$$

Forward algorithm

1. Base case:

$$\alpha_u(1) = a_{\text{START},u} \quad (36)$$

2. Recursive case:

$$\alpha_u(i) = \sum_v \alpha_v(i-1) \cdot a_{v,u} \cdot b_u(x_i) \cdot c_{x_i}(y_i) \quad (37)$$

Backward algorithm

1. Base case:

$$\beta_u(n) = a_{u,\text{STOP}} \cdot b_u(x_n) \cdot c_{x_n}(y_n) \quad (38)$$

2. Recursive case:

$$\beta_u(i) = \sum_v \beta_v(i+1) \cdot a_{v,u} \cdot b_u(x_i) \cdot c_{x_i}(y_i) \quad (39)$$

For this algorithm, time complexity will be $O(n|T|^2)$ for both forward and backward algorithms. We need to visit each node in the layer once, giving $|T|^2$ total operations. We also traverse from **START** to **STOP**, over n total layers. Hence the total number of operations is $n \cdot |T|^2$.