

Lecture 2: Hidden Markov Models

A hidden Markov model (HMM) consists of a set of internal states and a set of observable tokens. A run of a hidden Markov model generates a hidden state sequence s_1, \dots, s_T and a sequence of observable tokens a_1, \dots, a_T .

$$\pi(s) = P(s_1 = s), \sum_s \pi(s) = 1$$

$$T(w|s) = P(s_{t+1} = w | s_t = s), \sum_w T(w|s) = 1$$

$$O(a|s) = P(a_t = a | s_t = s), \sum_a O(a|s) = 1$$

$$P(s_1, \dots, s_T, a_1, \dots, a_T) = \pi(s_1) \left(\prod_{t=1}^T O(a_t | s_t) \right) \left(\prod_{t=1}^{T-1} T(s_{t+1} | s_t) \right)$$

Applications of HMMs:

- Speech Recognition. The hidden states are word positions and the observable tokens are acoustic feature vectors.

- Part of speech tagging. The hidden states are the parts of speech (noun, verb, adjective, and so on).

- DNA sequence analysis. The hidden states might be protein secondary structure or a position in a homologous sequence.

1 The Viterbi Algorithm

$$\text{Viterbi}[s, t] = \max_{s_1, \dots, s_{t-1}} P(s_1, \dots, s_{t-1}, s, a_1, \dots, a_{t-1})$$

$$\text{Viterbi}[s, 1] = \pi(s)$$

$$\text{Viterbi}[w, t + 1] = \max_s \text{Viterbi}[s, t] O(a_t | s) T(w | s)$$

$$s_T^* = \operatorname{argmax}_s \text{Viterbi}(s, T) O(a_T | s)$$

The best predecessor s_t can be recorded for each possible value of s_{t+1} and the best path can be constructed by working backward from s_T^* through best predecessors.

2 The Forward-Backward Procedure

$$\text{Forward}[s, t] = P(a_1, \dots, a_{t-1}, s_t = s)$$

$$\text{Backward}[s, t] = P(a_t, \dots, a_T | s_t = s)$$

$$\text{Forward}[s, 1] = \pi(s)$$

$$\text{Forward}[w, t + 1] = \sum_s \text{forward}[s, t] O(a_t | s) T(w | s)$$

$$\text{Backward}[s, T] = O(a_T | s)$$

$$\text{Backward}[s, t] = O(a_t | s) \sum_w T(w | s) \text{Backward}[s, t + 1]$$

$$P(a_1, \dots, a_T) = \sum_s \pi(s) \text{Backward}[s, 1]$$

$$P(s_t = s | a_1, \dots, a_T) = \frac{\text{Forward}(s, t) \text{Backward}[s, t]}{P(a_1, \dots, a_T)}$$

3 Trigram Language Models

Let $\#(w)$ be the number of times that the word w appears in a certain training corpus. Let $\#(w_1w_2)$ be the number of times that the pair of words w_1w_2 occurs and similarly for $\#(w_1, w_2, w_3)$ for the triple of words w_1, w_2, w_3 . Let N be the total number of word occurrences. A interpolated trigram model predicts the word w_3 following a given pair w_1, w_2 as follows.

$$P(w_3|w_1, w_2) = \lambda_1 \left(\frac{\#(w_1, w_2, w_3)}{\#(w_1, w_2)} \right) + \lambda_2 \left(\frac{\#(, w_2, w_3)}{\#(w_2)} \right) + \lambda_3 \left(\frac{\#(w_3)}{N} \right) \quad (1)$$

Here λ_1, λ_2 , and λ_3 are non-negative weights which sum to one:

$$\lambda_1 + \lambda_2 + \lambda_3 = 1$$

A weighted sum, such as (??), where the weights are non-negative and sum to one, is called a *convex combination*. Any convex combination of probability distributions is also a probability distribution. A convex combination of distributions is often called an *interpolated model*. In a trigram language model the weights λ_1, λ_2 and λ_3 are usually taken to depend on some way on the pair w_1, w_2 . This is ok since we can hold w_1, w_2 fixed in defining the conditional distribution $P(w_3|w_1, w_2)$.

A trigram language model defines the hidden states used in standard HMM-based speech recognition. A hidden state is a triple of words w_1, w_2, w_3 together with an index for a position in the last word. For example “we the pe*ople” is a state specifying that the two preceding words were “we” and “the” and that we are currently at the o in the word “people” so we should expect to be hearing an “o” sound. The state transition probabilities can be taken to be the following.

$$T(w_1, w_2, [\alpha_1 \dots \alpha_i * \alpha_{i+1} \alpha_{i+2} \dots \alpha_k] \mid w_1, w_2, [\alpha_1 \dots \alpha_i * \alpha_{i+1} \alpha_{i+2} \dots \alpha_k]) = 1/2$$

$$T(w_1, w_2, [\alpha_1 \dots \alpha_i \alpha_{i+1} * \alpha_{i+2} \dots \alpha_k] \mid w_1, w_2, [\alpha_1 \dots \alpha_i * \alpha_{i+1} \dots \alpha_k]) = 1/2$$

$$T(w_2, w_3, *w_4 \mid w_1, w_2, w_3*) = P(w_4|w_2, w_3)$$

The output probabilities are determined by a “acoustic model” specifying the probability distribution over acoustic feature vectors given the current phoneme of the hidden state. Actually, the distribution on acoustic features is usually taken to depend on a “triphone” — the preceding phoneme, the current phoneme, and the next phoneme. Pentaphones are even used in some systems.

4 Problem

Suppose that we have two hidden states A and B and two observable symbols a and b and an HMM defined by the following probabilities.

$$\pi(A) = \pi(B) = 1/2$$

$$T(A|A) = T(B|B) = 1 - \epsilon$$

$$T(B|A) = T(A|B) = \epsilon$$

$$O(a|A) = O(b|B) = 1 - \delta$$

$$O(b|A) = O(a|B) = \delta$$

Now suppose that we observe a sequence of T a 's. Let $F(A, t)$ abbreviate Forward(A, t) and similarly for $F(B, t)$. Give the values of $F(A, 1)$ and $F(B, 1)$ and give equation for computing $F(A, t+1)$ and $F(B, t+1)$ as a function of $F(A, t)$ and $F(B, t)$. Similarly let $B(A, t)$ abbreviate Backward(A, t). Give the values of $B(A, T)$ and $B(B, T)$ and give equation for computing $B(A, t)$ and $B(B, t)$ as a function of $F(A, t+1)$ and $F(B, t+1)$.

Extra Credit: Solve for $F(A, t)$, $B(A, t)$ and $P(a_1, \dots, a_T)$. Graph $P(s_t = A)$ as a function of t for $\epsilon = \delta = 1/4$, and $T = 20$ (you should not calculate the actual numbers for $P(s_t = A)$ if you can see qualitatively what the graph must look like).