

CSE 150. Assignment 6

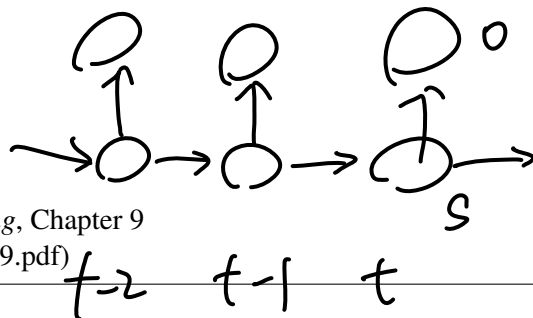
Spring 2018

Out: Tue May 22

Due: Tue May 29

Supplementary reading:

- Russell & Norvig, Chapter 15.
- Jurafsky & Martin, *Speech and Language Processing*, Chapter 9
(available at <https://web.stanford.edu/jurafsky/slp3/9.pdf>)



6.1 Conditional independence

Consider the hidden Markov model (HMM) shown below, with hidden states S_t and observations O_t for times $t \in \{1, 2, \dots, T\}$. Indicate whether the following statements are true or false.

F
T
F
F
T
F
T
T
T
F
F
T

$$P(S_t | S_{t-1}) = P(S_t | S_{t-1}, O_t)$$

$$P(S_t | S_{t-1}) = P(S_t | S_{t-1}, O_{t-1})$$

$$P(S_t | S_{t-1}) = P(S_t | S_{t-1}, S_{t+1})$$

$$P(S_t | O_{t-1}) = P(S_t | O_1, O_2, \dots, O_{t-1})$$

$$P(O_t | S_{t-1}) = P(O_t | S_{t-1}, O_{t-1})$$

$$P(O_t | O_{t-1}) = P(O_t | O_1, O_2, \dots, O_{t-1})$$

$$P(O_1, O_2, \dots, O_T) = \prod_{t=1}^T P(O_t | O_1, \dots, O_{t-1}) \quad \text{Handwritten: } P(O_1)P(O_2|O_1)P(O_3|O_1, O_2)$$

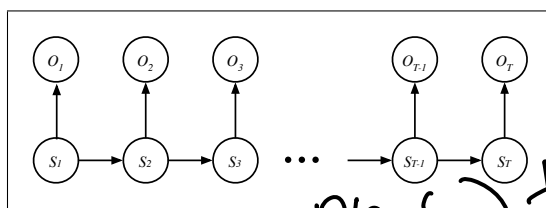
$$P(S_2, S_3, \dots, S_T | S_1) = \prod_{t=2}^T P(S_t | S_{t-1})$$

$$P(S_1, S_2, \dots, S_{T-1} | S_T) = \prod_{t=1}^{T-1} P(S_t | S_{t+1}) \quad \text{Handwritten: } P(S_1|S_2)P(S_2|S_3) \dots P(S_{T-1}|S_T)$$

$$P(S_1, S_2, \dots, S_T | O_1, O_2, \dots, O_T) = \prod_{t=1}^T P(S_t | O_t)$$

$$P(S_1, S_2, \dots, S_T, O_1, O_2, \dots, O_T) = \prod_{t=1}^T P(S_t, O_t)$$

$$P(O_1, O_2, \dots, O_T | S_1, S_2, \dots, S_T) = \prod_{t=1}^T P(O_t | S_t)$$



Handwritten notes:
 $P(O_1, S_1)$
 $P(O_2, S_2 | S_1, O_1)$
 $P(O_3, S_3 | S_1, O_1, O_2, S_2)$
 $P(O_T, S_T | S_1, O_1, \dots, O_{T-1}, S_{T-1})$

6.2 More conditional independence

Indicate the **smallest** subset of evidence nodes that must be considered to compute each conditional probability shown below. The first two problems are done as examples. (You may assume everywhere that $2 < t < T - 1$: i.e., do not worry about special boundary cases.)

$$P(S_t | S_1, S_2, \dots, S_{t-1}) = P(S_t | S_{t-1})$$

$$P(O_t | S_1, S_2, \dots, S_T) = P(O_t | S_t)$$

$$P(S_t | S_{t+1}, S_{t+2}, \dots, S_T) = \frac{P(S_t | S_{t+1})}{P(S_t | S_{t+1})}$$

$$P(S_t | O_t, O_{t-1}, O_{t+1}) = \frac{P(S_t | O_t, O_{t-1}, O_{t+1})}{P(S_t | O_t, O_{t-1}, O_{t+1})}$$

$$P(S_t | O_t, O_{t-1}, O_{t+1}, S_{t-1}, S_{t+1}) = \frac{P(S_t | S_{t-1}, S_{t+1}, O_t)}{P(S_t | S_{t-1}, S_{t+1}, O_t)}$$

$$P(S_t | S_1, S_T, O_1, O_t, O_T) = \frac{P(S_t | S_1, S_T)}{P(S_t | S_1, S_T)}$$

$$P(S_t | O_t, O_{t+1}, \dots, O_T) = \frac{P(S_t | O_t, O_{t+1}, \dots, O_T)}{P(S_t | O_t, O_{t+1}, \dots, O_T)}$$

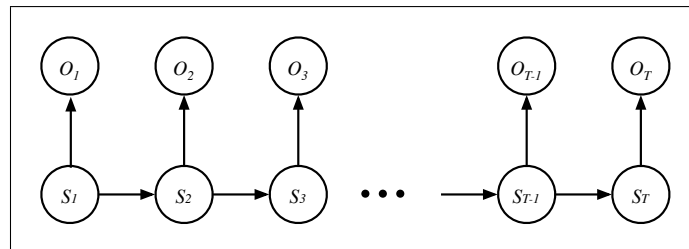
$$P(O_t | O_1, O_2, \dots, O_{t-1}) = \frac{P(O_t | O_1, O_2, \dots, O_{t-1})}{P(O_t | O_1, O_2, \dots, O_{t-1})}$$

$$P(O_t | O_1, O_2, \dots, O_{t-1}, S_{t-1}) = \frac{P(O_t | S_{t-1})}{P(O_t | S_{t-1})}$$

$$P(O_t | O_1, O_2, \dots, O_{t-1}, S_{t-2}) = \frac{P(O_t | S_{t-2}, O_{t-1})}{P(O_t | S_{t-2}, O_{t-1})}$$

$$P(O_t | S_{t-2}, S_{t-1}, S_{t+1}, S_{t+2}) = \frac{P(O_t | S_{t-1}, S_{t+1})}{P(O_t | S_{t-1}, S_{t+1})}$$

$$P(O_t | O_{t-1}, O_{t+1}, S_1, S_T) = \frac{P(O_t | O_{t-1}, O_{t+1}, S_1, S_T)}{P(O_t | O_{t-1}, O_{t+1}, S_1, S_T)}$$



6.3 Viterbi algorithm

In this problem, you will decode an English phrase from a long sequence of non-text observations. To do so, you will implement the same algorithm used in modern engines for automatic speech recognition. In a speech recognizer, these observations would be derived from real-valued measurements of acoustic waveforms. Here, for simplicity, the observations only take on binary values, but the high-level concepts are the same.

Consider a discrete HMM with $n = 27$ hidden states $S_t \in \{1, 2, \dots, 27\}$ and binary observations $O_t \in \{0, 1\}$. Download the ASCII data files from the course web site for this assignment. These files contain parameter values for the initial state distribution $\pi_i = P(S_1 = i)$, the transition matrix $a_{ij} = P(S_{t+1} = j | S_t = i)$, and the emission matrix $b_{ik} = P(O_t = k | S_t = i)$, as well as a long bit sequence of $T = 308000$ observations.

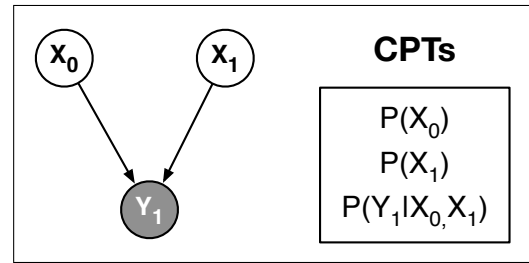
Use the Viterbi algorithm to compute the most probable sequence of hidden states conditioned on this particular sequence of observations. As always, you may program in the language of your choice. Turn in the following:

- (a) **a hard-copy print-out of your source code**
- (b) **a plot of the most likely sequence of hidden states versus time.**

To check your answer: suppose that the hidden states $\{1, 2, \dots, 26\}$ represent the letters $\{a, b, \dots, z\}$ of the English alphabet, and suppose that hidden state 27 encodes a space between words. If you have implemented the Viterbi algorithm correctly, the most probable sequence of hidden states (*ignoring repeated elements*) will reveal a highly recognizable message, as well as an interesting commentary on our times.

6.4 Belief updating

Consider the simple belief network on the right with nodes X_0 , X_1 , and Y_1 . To compute the posterior probability $P(X_1|Y_1)$, we can use Bayes rule:



(a) Show how to compute the term $P(Y_1|X_1)$ in the numerator of Bayes rule.

(b) Show how to compute the term $P(Y_1)$ in the denominator of Bayes rule.

Now consider the belief network shown at the bottom of the page. It does not have the same structure as an HMM, but using similar ideas we can derive efficient algorithms for inference. In particular, consider how to compute the posterior probability $P(X_t|Y_1, Y_2, \dots, Y_t)$ that accounts for evidence up to and including time t . We can derive an efficient recursion from Bayes rule:

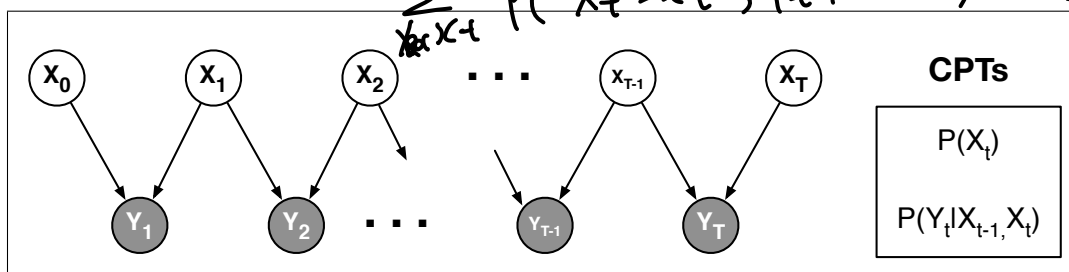
$$P(X_t|Y_1, Y_2, \dots, Y_t) = \frac{P(Y_t|X_t, Y_1, Y_2, \dots, Y_{t-1}) P(X_t|Y_1, Y_2, \dots, Y_{t-1})}{P(Y_t|Y_1, \dots, Y_{t-1})}$$

where the nodes Y_1, Y_2, \dots, Y_{t-1} are treated as background evidence. In parts (c-e) of this problem you will compute the individual terms that appear in this version of Bayes rule. You should express your answers in terms of the CPTs of the belief network and the probabilities $P(X_{t-1}=x|Y_1, Y_2, \dots, Y_{t-1})$ which you may assume have been computed at a previous step of the recursion. Your answers to parts (a) and (b) may be instructive for parts (d) and (e).

(c) Show how to simplify the term $P(X_t|Y_1, Y_2, \dots, Y_{t-1})$ in the numerator of Bayes rule.

(d) Show how to compute the term $P(Y_t|X_t, Y_1, Y_2, \dots, Y_{t-1})$ in the numerator of Bayes rule.

(e) Show how to compute the term $P(Y_t|Y_1, Y_2, \dots, Y_{t-1})$ in the denominator of Bayes rule.



6.5 Inference in HMMs

Consider a discrete HMM with hidden states S_t , observations O_t , transition matrix $a_{ij} = P(S_{t+1} = j | S_t = i)$ and emission matrix $b_{ik} = P(O_t = k | S_t = i)$. In class, we defined the forward-backward probabilities:

$$\alpha_{it} = P(o_1, o_2, \dots, o_t, S_t = i),$$

$$\beta_{it} = P(o_{t+1}, o_{t+2}, \dots, o_T | S_t = i),$$

for a particular observation sequence $\{o_1, o_2, \dots, o_T\}$ of length T . In terms of these probabilities, which you may assume to be given, as well as the transition and emission matrices of the HMM, show how to (efficiently) compute the following quantities:

- $P(S_{t+1} = j | S_t = i, o_1, o_2, \dots, o_T)$
- $P(S_t = i | S_{t+1} = j, o_1, o_2, \dots, o_T)$
- $P(S_{t-1} = i, S_t = j, S_{t+1} = k | o_1, o_2, \dots, o_T)$
- $\hat{s}_t = \operatorname{argmax}_i \left[P(S_t = i | o_1, o_2, \dots, o_T) \right]$

Handwritten derivation for (a):

$$= \frac{P(S_{t+1} = j | S_t = i, o_{t+1}, \dots, o_T)}{P(o_{t+1}, \dots, o_T | S_t = i)} \cdot \frac{P(o_{t+1}, \dots, o_T | S_t = i)}{P(o_{t+1}, \dots, o_T | S_t = i)}$$

Handwritten derivation for (b):

$$= \frac{P(S_t = i | S_{t+1} = j, o_1, o_2, \dots, o_T)}{P(o_1, o_2, \dots, o_T | S_t = i)} \cdot \frac{P(o_1, o_2, \dots, o_T | S_t = i)}{P(o_1, o_2, \dots, o_T | S_t = i)}$$

Handwritten derivation for (d):

$$\hat{s}_t = \operatorname{argmax}_i \left[\frac{\alpha_{it} \beta_{it}}{P(o_1, o_2, \dots, o_T)} \right]$$

In all these problems, you may assume that $t > 1$ and $t < T$; in particular, you are *not* asked to consider the boundary cases.

6.6 Most likely hidden states

The Viterbi algorithm in HMMs computes the most likely *sequence* of hidden states for a particular sequence of observations:

$$\{s_1^*, s_2^*, \dots, s_T^*\} = \operatorname{argmax}_{\{s_1, s_2, \dots, s_T\}} \left[P(s_1, s_2, \dots, s_T | o_1, o_2, \dots, o_T) \right]$$

Consider how these *collectively* optimal hidden states s_t^* differ (if at all) from the *individually* optimal hidden states \hat{s}_t in part (d) of the previous problem:

$$\hat{s}_t = \operatorname{argmax}_i \left[P(S_t = i | o_1, o_2, \dots, o_T) \right].$$

Answer the following [yes/no] questions:

- Is it possible that $P(\hat{s}_1, \dots, \hat{s}_T | o_1, \dots, o_T) > P(s_1^*, \dots, s_T^* | o_1, \dots, o_T)$?
- Is it possible that $\hat{s}_t = s_t^*$ for all t ?
- Is it necessarily true that $\hat{s}_t = s_t^*$ for all t ?
- Is it necessarily true that $P(\hat{s}_1, \hat{s}_2, \dots, \hat{s}_T) > 0$?