

COMP 5212 Machine Learning (2024 Fall)

Homework 3:

Hand out: November 8, 2024

Due: November 21, 2024, 11:59 PM

Total Points: 77

Tran Chun Yui

20963715

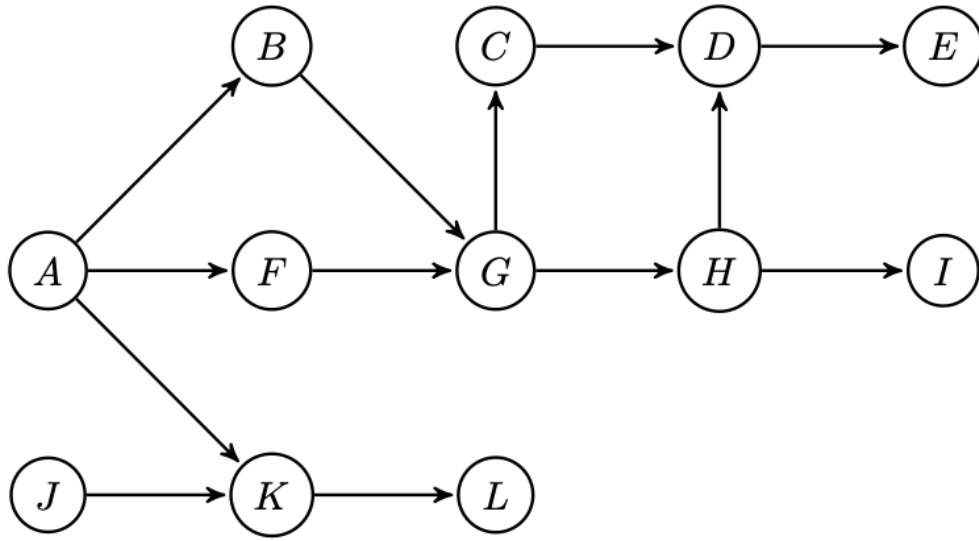
Your solution should contain below information at the top of its first page.

1. Your name
2. Your student id number

Some Notes:

- Homeworks will not be easy, please start early.
- For this homework, please submit a single PDF file to Canvas. Make sure to prepare your solution to each problem on a separate page.
- The total points of this homework is 100, and a score of 100 already gives you full grades on this homework. There are possibly bonus questions, you can optionally work on them if you are interested and have time.
- You can choose either using \LaTeX by inserting your solutions to the problem pdf, or manually write your solutions on clean white papers and scan it as a pdf file – in the case of handwriting, please write clearly and briefly, we are not responsible to extract information from unclear handwriting. **We highly recommend you use \LaTeX for the sake of any misunderstandings about the handwriting.** If your submission is a scan of a handwritten solution, make sure that it is of high enough resolution to be easily read. At least 300dpi and possibly denser.
- We encourage students to work in groups for homeworks, but the students need to write down the homework solutions or the code independently. In case that you work with others on the homework, please write down the names of people with whom you've discussed the homework. You are not allowed to copy, refer to, or look at the exact solutions from previous years, online, or other resources.
- **Late Policy:** 3 free late days in total across the semester, for additional late days, 20% penalization applied for each day late. **No assignment will be accepted more than 3 days late.**
- Please refer to the Course Logistics page for more details on the honor code and logisitcs. **We have zero tolerance — in the case of honor code violation for a single time, you will fail this course directly.**

1 Probabilistic Graphical Models (24 pts)



Which of the following statements are true given the network above? For true statements, brief explain why. For false statements, show one active trail using the Bayes Ball Algorithm.

1. $P(H, J) = P(H)P(J)$

False. Bayes Ball can be passed from J-K-A-F-G-H so H and J are not independent.

2. $P(H, J|L) = P(H|L)P(J|L)$

False. Bayes Ball can still be passed from J-K-A-F-G-H so H and J are not conditionally independent given L.

3. $P(C, I|F) = P(C|F)P(I|F)$

False. Bayes Ball can be passed from C-D-H-I so it is not conditionally independent given F.

4. $P(C, I|G, E) = P(C|G, E)P(I|G, E)$

False. Bayes Ball can still be passed from C-D-H-I so it is not conditionally independent given G and E.

5. $P(A, D|B) = P(A|B)P(D|B)$

False. Bayes Ball can be passed from A-F-G-C-D so it is not conditionally independent given B.

6. $P(B, F) = P(B)P(F)$

False. Bayes Ball can be passed from B-G-F so they are not independent.

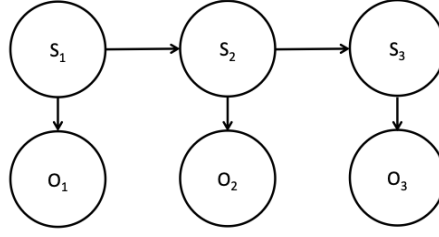
7. $P(C, K|B, F) = P(C|B, F)P(K|B, F)$

True. B and F have fully explained G so K is not affecting C.

8. $P(E, K|L) = P(E|L)P(K|L)$

False. Bayes Ball can be passed from K-A-F-G-H-D-E so they are not conditionally independent given L.

2 Hidden Markov Models (53 pts)



1. [5 points] Assume that we have the Hidden Markov Model (HMM) depicted in the figure above. If each of the states can take on k different values and a total of m different observations are possible for each state, how many parameters are required to fully define this HMM? Justify your answer.

HMM shares parameters for different t . The parameters needed are :

1. initial probability π_i
2. k by k transition matrix
3. k by m emission matrix.

Since the matrix are row stochastic, one column can be inferred by others columns. Total parameters needed are $(k - 1) + k(k - 1) + k(m - 1) = k^2 + km - k - 1$

2. [10 points] In class, we have learned the forward algorithm to compute the probability of observations, $P(\{O_t\}_{t=1}^T)$. Please try to derive a different algorithm to compute $P(\{O_t\}_{t=1}^T)$ in the backward direction. Similar to the forward algorithm in the slides, you should indicate the initialization, the computation equations in the backward direction at each step, and the termination. (While the figure above only shows a sequence length of 3, please derive a general form here assuming the sequence length is T)

$$2. P(\{O_t\}_{t=1}^T) = \sum_{S_t} P(\{O_t\}_{t=1}^T, S_t) = \sum_{S_t} P(O_1, \dots, O_t, S_t) P(O_{t+1}, \dots, O_T, S_t)$$

Let $t = 1$:

$$\begin{aligned}
 &= \sum_k P(O_1, S_1 = k) P(O_2, \dots, O_T, S_1 = k) \\
 &= \sum_k P(O_1 | S_1 = k) P(S_1 = k) P(O_2, \dots, O_T, S_1 = k)
 \end{aligned}$$

Let $\beta_t^k = P(O_{t+1}, \dots, O_T, S_t = k)$

$$P(\{O_t\}_{t=1}^T) = \sum_k P(O_1 | S_1 = k) P(S_1 = k) \beta_1^k$$

Similar to the decoding backward algorithm:

initialize (for all k):

$$\beta_T^k = 1$$

iterate for $t = T-1$ to 1

$$\beta_t^k = \sum_i P(S_{t+1} = i | S_t = k) P(O_{t+1} | S_{t+1} = i) \beta_{t+1}^i$$

Termination:

$$P(\{O_t\}_{t=1}^T) = \sum_k P(O_1 | S_1 = k) P(S_1 = k) \beta_1^k$$

3. **[10 points]** Similar to the last question, in class, we have learned the forward-backward algorithm to compute the posterior distribution $P(S_t = k | \{O_t\}_{t=1}^T)$. Please try to derive a backward-forward version of it – where you eliminate the variables on the left side of S_t in the backward direction, and the right side of S_t in the forward direction. Similar to that in the slides, you should indicate the initialization, the computation equations, and the termination. (While the figure above only shows a sequence length of 3, please derive a general form here assuming the sequence length is T)

3.

$$P(S_t | \{O_t\}_{t=1}^T) = \frac{P(S_t, \{O_t\}_{t=1}^T)}{P(\{O_t\}_{t=1}^T)}$$

where $P(S_t, \{O_t\}_{t=1}^T) = P(O_1, \dots, O_t, S_t) P(O_{t+1}, \dots, O_T | S_t)$

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

$$\begin{aligned} P(O_{t+1}, \dots, O_T | S_t) &= \sum_{S_{t+1} \dots S_T} P(O_{t+1}, \dots, O_T, S_{t+1}, \dots, S_T | S_t) \\ &= \sum_{S_{t+1} \dots S_T} \prod_i P(O_i | S_i) \prod_{i=t+1} P(S_{i+1} | S_i) \end{aligned}$$

by doing marginalization from $t+1$ to T :

$$= \sum_T P(O_T | S_T) \sum_{T-1} P(O_{T-1} | S_{T-1}) P(S_T | S_{T-1}) \dots \sum_{t+1} P(O_{t+1} | S_{t+1}) P(S_{t+1} | S_t) P(S_{t+2} | S_{t+1})$$

Let $\beta_{t+1}^i = P(S_{t+1} = i | S_t = k)$ The forward algorithm is defined as:

Initialize: $\beta_{t+1} = P(S_{t+1} = i | S_t = k)$ if O_{t+1} is not observed, set to 1

Update: $\beta_j^i = \sum_k P(O_{j-1} | S_{j-1} = k) P(S_j | S_{j-1} = k) \beta_{j-1}^k$

Terminate: $P(O_{t+1}, \dots, O_T | S_t = k) = \beta_T^k = \sum_i P(O_T | S_T = k) \beta_{T-1}^i$

Now, for the backward algorithm, Consider $P(O_1, \dots, O_t, S_t)$:

$$\begin{aligned} P(O_1, \dots, O_t, S_t) &= \sum_{S_1 \dots S_{t-1}} P(O_1, \dots, O_t, S_1, \dots, S_t) \\ &= \sum_{S_1 \dots S_{t-1}} \prod_i P(O_i | S_i) \prod_{i=t-1} P(S_{i+1} | S_i) \end{aligned}$$

by doing marginalization from t-1 to 1:

$$= \sum_1 P(O_1|S_1) \sum_2 P(O_2|S_2)P(S_2|S_1) \dots \sum_{t-1} P(O_{t-1}|S_{t-1})P(S_{t-1}|S_{t-2})P(S_t|S_{t-1})P(O_t|S_t)$$

Initialize: $\alpha_{t-1}^i = P(S_t = k|S_{t-1} = i)P(O_t|S_t = k)$ if O_{t-1} is not observed, set to 1

Update: $\alpha_j^i = \sum_i P(O_j|S_j = i)P(S_j = i|S_{j-1} = k)\alpha_{j+1}^i$

Terminate: $P(O_1, \dots, O_t, S_t = k) = \alpha_1^k = \sum_i P(O_1|S_1 = i)\alpha_2^i$

After finding $P(O_1, \dots, O_t, S_t = k)$ and $P(O_{t+1}, \dots, O_T|S_t = k)$,

compute $P(S_t = k|\{O_t\}_{t=1}^T) = \frac{P(S_t=k, \{O_t\}_{t=1}^T)}{P(\{O_t\}_{t=1}^T)} = \frac{\alpha_1^k \beta_T^k}{\sum_i \alpha_1^i \beta_T^i}$

Suppose that we have binary states (labeled A and B) and binary observations (labeled 0 and 1) and the initial, transition, and emission probabilities are as given in the table in Figure 1.

| State | $P(S_1)$ |
|-------|----------|
| A | 0.99 |
| B | 0.01 |

(a) Initial probs.

| S_1 | S_2 | $P(S_2 S_1)$ |
|-------|-------|--------------|
| A | A | 0.99 |
| A | B | 0.01 |
| B | A | 0.01 |
| B | B | 0.99 |

(b) Transition probs.

| S | O | $P(O S)$ |
|-----|-----|----------|
| A | 0 | 0.8 |
| A | 1 | 0.2 |
| B | 0 | 0.1 |
| B | 1 | 0.9 |

(c) Emission probs.

Figure 1: Probabilities Table

4. [7 points] Using the forward algorithm, compute the probability that we observe the sequence $O_1 = 0, O_2 = 1$, and $O_3 = 0$. Show your work (i.e., show each of your alphas, you can directly use the equations from the lecture).

4. The initial forward variables are:

For $\alpha_0(1)$ (state 1):

$$\alpha_0(1) = \pi_1 \cdot B_{1,O_1} = 0.99 \cdot 0.8 = 0.792$$

For $\alpha_0(2)$ (state 2):

$$\alpha_0(2) = \pi_2 \cdot B_{2,O_1} = 0.01 \cdot 0.1 = 0.001$$

$$\alpha_1(1) = [0.792 \cdot 0.99 + 0.001 \cdot 0.01] \cdot 0.2 = [0.78308 + 0.00001] \cdot 0.2 = 0.156818$$

$$\alpha_1(2) = [0.792 \cdot 0.01 + 0.001 \cdot 0.99] \cdot 0.9 = [0.00792 + 0.00099] \cdot 0.9 = 0.008019$$

$$\alpha_2(1) = [0.156818 \cdot 0.99 + 0.008019 \cdot 0.01] \cdot 0.8 = [0.155249 + 0.00008019] \cdot 0.8 = 0.12426401$$

$$\alpha_2(2) = [0.156818 \cdot 0.01 + 0.008019 \cdot 0.99] \cdot 0.2 = [0.00156818 + 0.00795981] \cdot 0.2 = 0.0009507$$

$$P(O_1, O_2, O_3) = \sum_{i=1}^2 \alpha_2(i)$$

$$P(O_1, O_2, O_3) = 0.12426401 + 0.0009507 = 0.125214707$$

5. [7 points] Using the backward algorithm you derived in Question 2 to compute the probability that we observe the aforementioned sequence ($O_1 = 0$, $O_2 = 1$, and $O_3 = 0$). You can directly use the equation you just derived.

5.

$$\beta_3^A = \beta_3^B = 1$$

The recursive equation for the backward algorithm is:

$$\beta_t^k = \sum_i P(S_{t+1} = i \mid S_t = k) P(O_{t+1} \mid S_{t+1} = i) \beta_{t+1}^i$$

Beta at time 1:

$$\beta_1^A = 0.793, \quad \beta_1^B = 0.107$$

Beta at time 0:

$$\beta_0^A = 0.157977, \quad \beta_0^B = 0.096923$$

Now, using the backward algorithm:

$$P(\{O_t\}_{t=1}^T) = \sum_k P(O_1 \mid S_1 = k) P(S_1 = k) \beta_1^k$$

$$P(\{O_t\}_{t=1}^T) = 0.125214707$$

The α values are:(from t=1 to t=3)

$$\alpha = \begin{bmatrix} 0.792 & 0.001 \\ 0.156818 & 0.008019 \\ 0.12426401 & 0.0009507 \end{bmatrix}$$

The β values are:

$$\beta = \begin{bmatrix} 0.157977 & 0.096923 \\ 0.793 & 0.107 \\ 1. & 1. \end{bmatrix}$$

The Posterior Probabilities are:

$$\text{Posterior} = \begin{bmatrix} 9.99225946 \times 10^{-1} & 7.74054441 \times 10^{-4} \\ 9.93147506 \times 10^{-1} & 6.85249377 \times 10^{-3} \\ 9.92407449 \times 10^{-1} & 7.59255061 \times 10^{-3} \end{bmatrix}$$

The Most Likely States are:

$$\text{Most Likely States} = \begin{bmatrix} A \\ A \\ A \end{bmatrix}$$

7. **[9 points]** Use the Viterbi algorithm to compute (and report) the most likely sequence of states. Show your work (i.e., show each of your Vs).

7. Using the algorithm,

$$\text{Viterbi at time 0: } V_0 = [0.792, 0.001]$$

$$\text{Viterbi at time 1: } V_1 = [0.156816, 0.007128]$$

$$\text{Viterbi at time 2: } V_2 = [0.12419827, 0.00070567]$$

The most likely final state is A , since $V_2(A) = 0.12419827$ is greater than $V_2(B) = 0.00070567$.

$$\text{Most likely sequence of states: } [A, A, A]$$

The probability of the most likely sequence is the value of $V_2(A)$:

$$P(\text{most likely sequence}) = V_2(A) = 0.12419827$$

8. **[2 points]** Is the most likely sequence of states the same as the sequence comprised of the most likely setting for each individual state? Provide a 1-2 sentence justification for your answer.

8. No, the most likely sequence of states is not necessarily the same as the sequence comprised of the most likely setting for each individual state. This is because the Viterbi algorithm considers the entire sequence of observations and transitions between states, while the backward-forward only considers the most likely state at each time step independently.