

RECITATION 8

HIDDEN MARKOV MODELS AND BAYES NET

10-601: INTRODUCTION TO MACHINE LEARNING

11/3/2021

1 HMMs

You are given the following training data:

win_C league_C Liverpool_D

win_C Liverpool_D league_C

Liverpool_D win_C

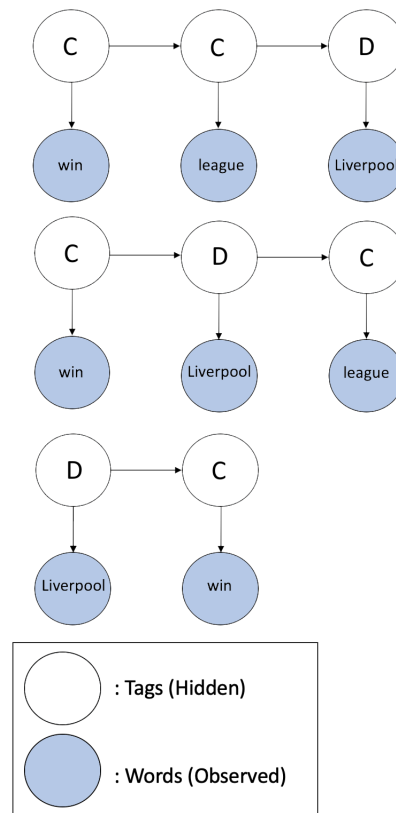


Figure 1: Visualization of Sequences

You are also given the following observed (validation) data: **Liverpool win league**

In this question, let each observed state $x_t \in \{1, 2, 3\}$, where 1 corresponds to **win**, 2 corresponds to **league**, and 3 corresponds to **Liverpool**. Let each hidden state $Y_t \in \{C, D\}$, where $s_1 = C$ and $s_2 = D$.

1. First, we need to train our HMM by generating the initial probabilities: $\boldsymbol{\pi}$, the transition probability matrix: \mathbf{B} , the emission probability matrix: \mathbf{A} .

(a) Find $\boldsymbol{\pi}$. Recall that $\pi_j = P(Y_1 = s_j)$.

- Find count matrix

$$\begin{array}{c} C \\ D \end{array} \begin{array}{c} \text{Count} \\ \left[\begin{array}{c} 2 \\ 1 \end{array} \right] \end{array} \xrightarrow{\text{Pseudocount}} \begin{array}{c} C \\ D \end{array} \begin{array}{c} \text{Count} \\ \left[\begin{array}{c} 3 \\ 2 \end{array} \right] \end{array}$$

- Get probability matrix $\boldsymbol{\pi}$:

$$\boldsymbol{\pi} = \begin{array}{c} C \\ D \end{array} \left[\begin{array}{c} 3/5 \\ 2/5 \end{array} \right]$$

(b) Find Transition Matrix: \mathbf{B} . Recall that $B_{jk} = P(Y_t = s_k \mid Y_{t-1} = s_j)$

- Find count matrix

$$\begin{array}{c} C \\ D \end{array} \begin{array}{cc} C & D \\ \left[\begin{array}{cc} 1 & 2 \\ 2 & 0 \end{array} \right] \end{array} \xrightarrow{\text{Pseudocount}} \begin{array}{c} C \\ D \end{array} \begin{array}{cc} C & D \\ \left[\begin{array}{cc} 2 & 3 \\ 3 & 1 \end{array} \right] \end{array}$$

- Get Transition Probability matrix \mathbf{B} :

$$\mathbf{B} = \begin{array}{c} C \\ D \end{array} \begin{array}{cc} C & D \\ \left[\begin{array}{cc} 2/5 & 3/5 \\ 3/4 & 1/4 \end{array} \right] \end{array}$$

(c) Find Emission Matrix: \mathbf{A} . Recall that $A_{jk} = P(X_t = k \mid Y_t = s_j)$.

- Find count matrix

$$\begin{array}{c} C \\ D \end{array} \begin{array}{ccc} \text{win} & \text{league} & \text{Liverpool} \\ \left[\begin{array}{ccc} 3 & 2 & 0 \\ 0 & 0 & 3 \end{array} \right] \end{array} \xrightarrow{\text{Pseudocount}} \begin{array}{c} C \\ D \end{array} \begin{array}{ccc} \text{win} & \text{league} & \text{Liverpool} \\ \left[\begin{array}{ccc} 4 & 3 & 1 \\ 1 & 1 & 4 \end{array} \right] \end{array}$$

- Get Emission Probability matrix \mathbf{A} :

$$\mathbf{A} = \begin{array}{c} C \\ D \end{array} \begin{array}{ccc} \text{win} & \text{league} & \text{Liverpool} \\ \left[\begin{array}{ccc} 1/2 & 3/8 & 1/8 \\ 1/6 & 1/6 & 2/3 \end{array} \right] \end{array}$$

2. What is the likelihood of observing this output?

Recall that:

$$\alpha_t(k) = P(x_{1:t}, Y_t = s_k)$$

$$\beta_t(k) = P(x_{t+1:T} | Y_t = s_k)$$

We also have the recursive procedure:

(a) $\alpha_1(j) = \pi_j A_{jx_1}$.

(b) For $t > 1$, $\alpha_t(j) = A_{jx_t} \sum_{k=1}^J \alpha_{t-1}(k) B_{kj}$

We want to find:

$$P(X_1 = \text{Liverpool}, X_2 = \text{win}, X_3 = \text{league})$$

$$= \sum_{y_t \in C, D} P(x_1 = \text{Liverpool}, x_2 = \text{win}, x_3 = \text{league}, Y_t = y_t)$$

$$= \sum_{y_t \in C, D} \alpha_3(y_t)$$

$$\alpha_1 = P(x_1, y_1) = P(x_1 | y_1) \cdot P(y_1) = A_{\cdot 3} \circ \pi$$

$$= \begin{bmatrix} \alpha_1(C) \\ \alpha_1(D) \end{bmatrix} = \begin{bmatrix} \pi_C * A_{C, x_1} \\ \pi_D * A_{C, x_1} \end{bmatrix} = \begin{bmatrix} \pi_C * A_{C, \text{Liverpool}} \\ \pi_D * A_{D, \text{Liverpool}} \end{bmatrix} = \begin{bmatrix} 1/8 \\ 2/3 \end{bmatrix} \circ \begin{bmatrix} 3/5 \\ 2/5 \end{bmatrix} = \begin{bmatrix} 0.075 \\ 0.26667 \end{bmatrix}$$

$$\alpha_2 = P(x_1, x_2, y_2) = P(x_2 | y_2) \cdot (P(y_2 | y_1) \cdot \alpha_1) = A_{\cdot 1} \circ (B^T \alpha_1)$$

$$= \begin{bmatrix} A_{C, \text{win}} * \sum_{y_t \in \{C, D\}} \alpha_1(y_t) * B_{y_t, C} \\ A_{D, \text{win}} * \sum_{y_t \in \{C, D\}} \alpha_1(y_t) * B_{y_t, D} \end{bmatrix} = \begin{bmatrix} 1/2 \\ 1/6 \end{bmatrix} \circ \left(\begin{bmatrix} 2/5 & 3/5 \\ 3/4 & 1/4 \end{bmatrix}^T \begin{bmatrix} 0.075 \\ 0.26667 \end{bmatrix} \right) = \begin{bmatrix} 0.11500125 \\ 0.01861125 \end{bmatrix}$$

$$\alpha_3 = A_{\cdot 2} \circ (B^T \alpha_2)$$

$$= \begin{bmatrix} 3/8 \\ 1/6 \end{bmatrix} \circ \left(\begin{bmatrix} 2/5 & 3/5 \\ 3/4 & 1/4 \end{bmatrix}^T \begin{bmatrix} 0.11500125 \\ 0.01861125 \end{bmatrix} \right) = \begin{bmatrix} 0.02248460156 \\ 0.01227559375 \end{bmatrix}$$

$$\text{Since } P(x_1 = \text{Liverpool}, x_2 = \text{win}, x_3 = \text{league}) = \sum_{y_t \in C, D} \alpha_3(y_t)$$

$$\therefore P(x_1 = \text{Liverpool}, x_2 = \text{win}, x_3 = \text{league})$$

$$= 0.02248460156 + 0.01227559375$$

$$= \boxed{0.03476019531}$$

You are now told that the observed data has the following tags:

Liverpool_D win_C league_D

3. Given the observed sequence of words (denote $\vec{x} = [\text{Liverpool}, \text{win}, \text{league}]^T$), what is the probability of these assigned tags $P(Y_1 = D|\vec{x})$, $P(Y_2 = C|\vec{x})$, $P(Y_3 = D|\vec{x})$?

Recall that:

$$P(Y_t = s_k|\vec{x}) = \frac{\alpha_t(s_k)\beta_t(s_k)}{P(\vec{x})}$$

So, we need to find β_T

We also have a similar recursive procedure

- (a) $\beta_T(j) = 1$ (All states could be ending states)
 (b) For $1 \leq t \leq T-1$, $\beta_t(j) = \sum_{k=1}^J A_{kx_{t+1}}\beta_{t+1}(k)B_{jk}$ (Generate x_{t+1} from any state)

Remember that: $\beta_t(s_k) = P(x_{t+1:T}|Y_t = s_k)$ and $\beta_T(s_k) = 1$

Using matrix notation:

$$\beta_2 = B(A_{,x_3} \circ \beta_3) = B(A_{,2} \circ \beta_3)$$

Recall that:

$$A_{,2} = \begin{bmatrix} 3/8 \\ 1/6 \end{bmatrix} \text{ and } \beta_3 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \text{ since } T = 3$$

$$\therefore \beta_2 = \begin{bmatrix} 0.25 \\ 0.3229 \end{bmatrix}$$

Now, we go on to solve β_1

$$\beta_1 = B(A_{,x_2} \circ \beta_2) = B(A_{,1} \circ \beta_2)$$

Again, recall that:

$$A_{,1} = \begin{bmatrix} 1/2 \\ 1/6 \end{bmatrix} \text{ and } \beta_2 = \begin{bmatrix} 0.25 \\ 0.3229 \end{bmatrix}$$

$$\therefore \beta_1 = \begin{bmatrix} 0.08229 \\ 0.1072 \end{bmatrix}$$

Now, we have our α and β matrix:

$$\alpha = \begin{array}{c} 1 \\ 2 \\ 3 \end{array} \begin{array}{cc} C & D \\ \left[\begin{array}{cc} 0.0750 & 0.26667 \\ 0.1150 & 0.0186 \\ 0.0225 & 0.0123 \end{array} \right] \end{array}$$

$$\beta = \begin{array}{c} 1 \\ 2 \\ 3 \end{array} \begin{array}{cc} C & D \\ \left[\begin{array}{cc} 0.0823 & 0.1072 \\ 0.2500 & 0.3229 \\ 1.0000 & 1.0000 \end{array} \right] \end{array}$$

$$\begin{aligned} P(Y_1 = D|\vec{x}) &= \frac{\alpha_1(D)\beta_1(D)}{P(\vec{x})} \\ &= \frac{0.26667 \times 0.1072}{0.03476019531} \\ &= 0.8224068865 \end{aligned}$$

$$\begin{aligned} P(Y_2 = C|\vec{x}) &= \frac{\alpha_2(C)\beta_2(C)}{P(\vec{x})} \\ &= \frac{0.1150 \times 0.2500}{0.03476019531} \\ &= 0.8270954678 \end{aligned}$$

$$\begin{aligned} P(Y_3 = C|\vec{x}) &= \frac{\alpha_3(C)\beta_3(C)}{P(\vec{x})} \\ &= \frac{0.0225 \times 1}{0.03476019531} \\ &= 0.6472921052 \end{aligned}$$

4. The sequence of words you observe is again the same: `Liverpool win league`
 However, you are only given the tag of the last word: `league_C`
 Using the Viterbi Algorithm, what is the most likely sequence of hidden states?

Recall that:

$$\omega_t(s_k) = \max_{y_{1:t-1}} P(x_{1:t}, y_{1:t-1}, y_t = s_k)$$

$$b_t(s_k) = \arg \max_{y_{1:t-1}} P(x_{1:t}, y_{1:t-1}, y_t = s_k)$$

Also, the recursive procedure for the Viterbi algorithm is as follows:

- (a) $\omega_0(s_k) = 1$ for $s_k = \text{START}$ and 0 for all other states.
- (b) For $t > 1$,
- $\omega_t(s_j) = \max_{1 \leq k \leq J} \omega_{t-1}(s_k) P(x_t | Y_t = s_j) P(Y_t = s_j | Y_{t-1} = s_k)$
 $= \max_{1 \leq k \leq J} \omega_{t-1}(s_k) A_{jx_t} B_{kj}$
 - $b_t(s_j) = \arg \max_{1 \leq k \leq J} \omega_{t-1}(s_k) P(x_t | Y_t = s_j) P(Y_t = s_j | Y_{t-1} = s_k)$
 $= \arg \max_{1 \leq k \leq J} \omega_{t-1}(s_k) A_{jx_t} B_{kj}$

What is the most likely sequence of tags given the observed data? (Select **C** if tie)

- (a) Set up the matrices ω and b

$$\omega = \begin{matrix} \omega_0 \\ \omega_1 \\ \omega_2 \\ \omega_3 \end{matrix} \begin{matrix} \text{C} & \text{D} & \text{START} \\ \left[\begin{array}{ccc} 0 & 0 & 1 \\ - & - & - \\ - & - & - \\ - & - & - \end{array} \right] \end{matrix}$$

and

$$b = \begin{matrix} b_1 \\ b_2 \\ b_3 \end{matrix} \begin{matrix} \text{C} & \text{D} \\ \left[\begin{array}{cc} - & - \\ - & - \\ - & - \end{array} \right] \end{matrix}$$

Initialize $w_0(\text{START}) = 1$

- (b) Solve for matrix entries using Dynamic Programming:

$$\begin{aligned}
 \omega_1(C) &= \max_{s_j \in \{C, D, \text{START}\}} P(x_1 = \text{Liverpool} | Y_1 = C) \omega_0(s_j) P(Y_1 = C) \\
 &= \frac{1}{8} \cdot 1 \cdot \frac{3}{5} \\
 &= \frac{3}{40}
 \end{aligned}$$

$$b_1(C) = \text{START}$$

$$\begin{aligned}
 \omega_1(D) &= \max_{s_j \in \{C, D, \text{START}\}} P(x_1 = \text{Liverpool} | Y_1 = D) \omega_0(s_j) P(Y_1 = D) \\
 &= \frac{2}{3} \cdot 1 \cdot \frac{2}{5} \\
 &= \frac{4}{15}
 \end{aligned}$$

$$b_1(D) = \text{START}$$

$$\begin{aligned}
 \omega_2(C) &= \max_{s_j \in \{C, D\}} P(x_2 = \text{win} | Y_2 = C) \omega_1(s_j) P(Y_2 = C | Y_1 = s_j) \\
 &= \max \left(\frac{1}{2} \cdot \frac{3}{40} \cdot \frac{2}{5}, \frac{1}{2} \cdot \frac{4}{15} \cdot \frac{3}{4} \right) \\
 &= \frac{1}{10}
 \end{aligned}$$

$$b_2(C) = D$$

$$\begin{aligned}
 \omega_2(D) &= \max_{s_j \in \{C, D\}} P(x_2 = \text{win} | Y_2 = D) \omega_1(s_j) P(Y_2 = D | Y_1 = s_j) \\
 &= \max \left(\frac{1}{6} \cdot \frac{3}{40} \cdot \frac{3}{5}, \frac{1}{6} \cdot \frac{4}{15} \cdot \frac{1}{4} \right) \\
 &= \frac{1}{90}
 \end{aligned}$$

$$b_2(D) = D$$

$$\begin{aligned}
 \omega_3(C) &= \max_{s_j \in \mathbf{C}, \mathbf{D}} P(x_3 = \text{league} | Y_3 = \mathbf{C}) \omega_2(s_j) P(Y_3 = \mathbf{C} | Y_2 = s_j) \\
 &= \max \left(\frac{3}{8} \cdot \frac{1}{10} \cdot \frac{2}{5}, \frac{3}{8} \cdot \frac{1}{90} \cdot \frac{3}{4} \right) \\
 &= \frac{3}{200}
 \end{aligned}$$

$$b_3(\mathbf{C}) = \mathbf{C}$$

$$\begin{aligned}
 \omega_3(D) &= \max_{s_j \in \mathbf{C}, \mathbf{D}} P(x_3 = \text{league} | Y_3 = \mathbf{D}) \omega_2(s_j) P(Y_3 = \mathbf{D} | Y_2 = s_j) \\
 &= \max \left(\frac{1}{6} \cdot \frac{1}{10} \cdot \frac{3}{5}, \frac{1}{6} \cdot \frac{1}{90} \cdot \frac{1}{4} \right) \\
 &= \frac{1}{100}
 \end{aligned}$$

$$b_3(\mathbf{D}) = \mathbf{C}$$

Now, to figure out the order, we set $\hat{y}_t = b_{t+1}(\hat{y}_{t+1})$

$$\begin{aligned}
 y_{T+1} &= \text{END} \\
 y_3 &= \mathbf{C} \\
 \hat{y}_2 &= b_3(\mathbf{C}) \\
 &= \mathbf{C} \\
 \hat{y}_1 &= b_2(\mathbf{C}) \\
 &= \mathbf{D} \\
 \hat{y}_0 &= b_1(\mathbf{D}) \\
 &= \text{START}
 \end{aligned}$$

So, the most likely sequence is **START-D-C-C-END**

2 Working in Log-space

2.1 Motivation

Given the following series of probability values:

$P(x_1 = 1)$	$P(x_2 = 1 \mid x_1 = 1)$	$P(x_3 = 1 \mid x_2 = 1, x_1 = 1)$
0.002	0.004	0.003

We want to find $P(x_1 = 1, x_2 = 1, x_3 = 1)$. Suppose we have a calculator which only has 4 decimal places of precision, so it can only store values of format X.XXXX

1. What is the correct value of $P(x_1 = 1, x_2 = 1, x_3 = 1)$ without any precision limits?

$$P(x_1 = 1, x_2 = 1, x_3 = 1) = P(x_3 = 1 \mid x_2 = 1, x_1 = 1) * P(x_2 = 1 \mid x_1 = 1) * P(x_1 = 1) \\ = 0.003 * 0.004 * 0.002 = 0.000000024$$

2. What is the value of $P(x_1 = 1, x_2 = 1, x_3 = 1)$ using our faulty calculator?

$$P(x_1 = 1, x_2 = 1) = P(x_2 = 1 \mid x_1 = 1)P(x_1 = 1) = 0.004 * 0.002 = 0.0000 \\ \implies \text{Truncated!}$$

$$P(x_1 = 1, x_2 = 1, x_3 = 1) = 0.0000 * 0.003 = 0.0000 \\ \implies \text{Truncated again!}$$

3. How do the values of $P(x_1 = 1, x_2 = 1, x_3 = 1)$ from part (1) and (2) compare?

$$\text{No precision limits: } P(x_1 = 1, x_2 = 1, x_3 = 1) = 0.000000024$$

$$\text{Faulty calculator: } P(x_1 = 1, x_2 = 1, x_3 = 1) = 0.0000$$

4. What is the value of $P(x_1 = 1, x_2 = 1, x_3 = 1)$ if we perform the same computation but in log space?

$$\log(P(x_1 = 1, x_2 = 1, x_3 = 1)) =$$

$$\log(P(x_1 = 1, x_2 = 1, x_3 = 1))$$

$$= \log(x_1 = 1) + \log(P(x_2 = 1 \mid x_1 = 1)) + \log(P(x_3 = 1 \mid x_2 = 1, x_1 = 1))$$

$$= \log(0.002) + \log(0.004) + \log(0.003)$$

$$= -6.2146 - 5.8091 - 5.5215$$

$$= -17.5452$$

$$\text{If we were to recover our value of } P(x_1 = 1, x_2 = 1, x_3 = 1) = e^{\log(P(x_1=1, x_2=1, x_3=1))} = \\ e^{-17.5452} = 0.000000024$$

This is good! But we can use the log sum exp trick to extend its use to even smaller scales.

2.2 Forward and Backward Algorithm in Log Space

In the forward algorithm, recall that the α 's can be computed using the recursive procedure:

- $\alpha_1(j) = \pi_j A_{jx_1}$
 - For $t > 1$, $\alpha_t(j) = A_{jx_t} \sum_{k=1}^J \alpha_{t-1}(k) B_{kj}$
1. Derive $\log(\alpha_1(j))$ in terms of $\log(\pi_j)$ and $\log(A_{jx_1})$

$$\log(\alpha_1(j)) = \log(\pi_j A_{jx_1}) = \log(\pi_j) + \log(A_{jx_1})$$
 2. Derive $\log(\alpha_t(j))$ in terms of $\log(\alpha_{t-1}(k))$ and $\log A_{kj}$

$$\begin{aligned} \log(\alpha_t(j)) &= \log\left(A_{jx_t} \sum_{k=1}^J \alpha_{t-1}(k) B_{kj}\right) \\ &= \log(A_{jx_t}) + \log\left(\sum_{k=1}^J \alpha_{t-1}(k) B_{kj}\right) \\ &= \log(A_{jx_t}) + \log\left(\sum_{k=1}^J e^{\log(\alpha_{t-1}(k) B_{kj})}\right) \\ &= \log(A_{jx_t}) + \log\left(\sum_{k=1}^J e^{\log(\alpha_{t-1}(k)) + \log(B_{kj})}\right) \end{aligned}$$

In the backward algorithm, we also have a similar recursive procedure:

- $\beta_T(j) = 1$
 - For $1 \leq t \leq T - 1$, $\beta_t(j) = \sum_{k=1}^J A_{kx_{t+1}} \beta_{t+1}(k) B_{jk}$
1. Derive $\log(\beta_T(j))$

$$\log(\beta_T(j)) = \log(1) = 0$$
 2. Derive $\log(\beta_t(j))$ in terms of $\log(A_{kx_{t+1}})$, $\log(\beta_{t+1}(k))$, and $\log(B_{jk})$

$$\begin{aligned} \log(\beta_t(j)) &= \log\left(\sum_{k=1}^J A_{kx_{t+1}} \beta_{t+1}(k) B_{jk}\right) \\ &= \log\left(\sum_{k=1}^J e^{\log(A_{kx_{t+1}} \beta_{t+1}(k) B_{jk})}\right) \\ &= \log\left(\sum_{k=1}^J e^{\log(A_{kx_{t+1}}) + \log(\beta_{t+1}(k)) + \log(B_{jk})}\right) \end{aligned}$$

3 Bayesian Networks

3.1 Practice problems

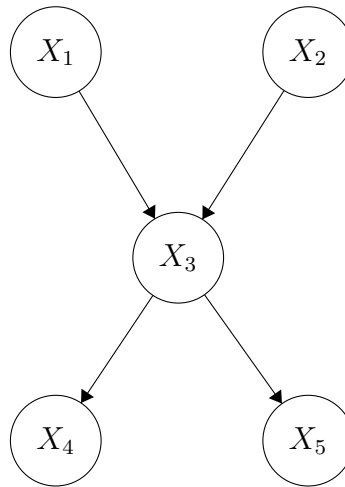


Figure 2: Graphical Model

1. Write down the factorization of the above directed graphical model.

$$\begin{aligned} &P(X_1, X_2, X_3, X_4, X_5) \\ &= P(X_1)P(X_2)P(X_3|X_1, X_2)P(X_4|X_3)P(X_5|X_3) \end{aligned}$$