

# RECITATION 8

## HIDDEN MARKOV MODELS

10-601: INTRODUCTION TO MACHINE LEARNING

4/1/2022

Version: 1.2

### 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

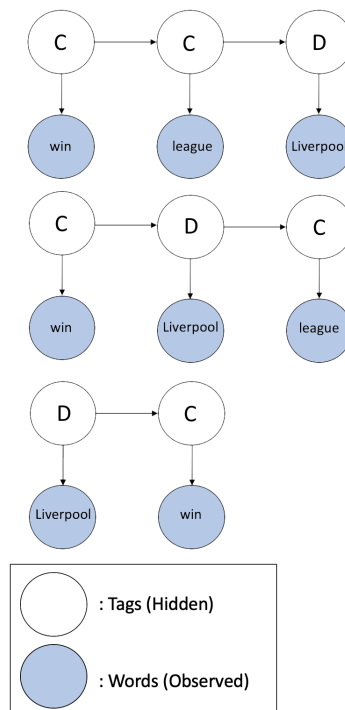


Figure 1: Visualization of Sequences

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

## 1.1 Initial, Emission, and Transition Matrices

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  $s_2 = D$ .

First, we need to estimate the HMM parameters - the initial probabilities:  $\pi$ , the transition probability matrix:  $\mathbf{B}$ , and the emission probability matrix:  $\mathbf{A}$ . Remember that we use MLE estimation to do so:

$$\hat{C}_k = \frac{\#(y_i^{(i)} = k)}{N} \quad \forall i, k$$

$$\hat{B}_{jk} = \frac{\#(y_t^{(i)} = k \text{ and } y_{t-1}^{(i)} = j)}{\#(y_{t-1}^{(i)} = j)} \quad \forall i, t > 1, j, k$$

$$\hat{A}_{jk} = \frac{\#(x_t^{(i)} = k \text{ and } y_t^{(i)} = j)}{\#(y_t^{(i)} = j)} \quad \forall i, t, j, k$$

*Note:* When learning an HMM, we add 1 to each count to make a pseudocount. This improves performance when evaluating unseen cases in the validation set or test set.

1. Find the initial matrix  $\pi$ . Recall that  $\pi_j = P(Y_1 = s_j)$ .

- Find count matrix and pseudocount matrix:

<i>Count</i>		<i>Count</i>
<i>C</i>	$\xrightarrow{\text{Pseudocount}}$	<i>C</i>
<i>D</i>		<i>D</i>

- Normalize:

$$\pi = \begin{matrix} C \\ D \end{matrix}$$

- Find count matrix and pseudocount matrix:

<i>Count</i>		<i>Count</i>
<i>C</i>	$\xrightarrow{\text{Pseudocount}}$	<i>C</i>
$\begin{bmatrix} 2 \\ 1 \end{bmatrix}$		$\begin{bmatrix} 3 \\ 2 \end{bmatrix}$
<i>D</i>		<i>D</i>

- Normalize:

$$\pi = \begin{matrix} C \\ D \end{matrix} \begin{bmatrix} 3/5 \\ 2/5 \end{bmatrix}$$

2. Find the transition matrix  $\mathbf{B}$ . Recall that  $B_{jk} = P(Y_t = s_k | Y_{t-1} = s_j)$ .

- Find count matrix and pseudocount matrix:

$$\begin{matrix} & C & D \\ C & & \\ D & & \end{matrix} \xrightarrow{\text{Pseudocount}} \begin{matrix} & C & D \\ C & & \\ D & & \end{matrix}$$

- Normalize:

$$\mathbf{B} = \begin{matrix} & C & D \\ C & & \\ D & & \end{matrix}$$

- Find count matrix and pseudocount matrix:

$$\begin{matrix} & C & D \\ C & \begin{bmatrix} 1 & 2 \end{bmatrix} \\ D & \begin{bmatrix} 2 & 0 \end{bmatrix} \end{matrix} \xrightarrow{\text{Pseudocount}} \begin{matrix} & C & D \\ C & \begin{bmatrix} 2 & 3 \end{bmatrix} \\ D & \begin{bmatrix} 3 & 1 \end{bmatrix} \end{matrix}$$

- Normalize:

$$B = \begin{matrix} & C & D \\ C & \begin{bmatrix} 2/5 & 3/5 \end{bmatrix} \\ D & \begin{bmatrix} 3/4 & 1/4 \end{bmatrix} \end{matrix}$$

3. Find the emission matrix  $\mathbf{A}$ . Recall that  $A_{jk} = P(X_t = k | Y_t = s_j)$ .

- Find count matrix and pseudocount matrix:

$$\begin{matrix} & \text{win} & \text{league} & \text{Liverpool} \\ C & & & \\ D & & & \end{matrix} \xrightarrow{\text{Pseudocount}} \begin{matrix} & \text{win} & \text{league} & \text{Liverpool} \\ C & & & \\ D & & & \end{matrix}$$

- Normalize:

$$\mathbf{A} = \begin{matrix} & \text{win} & \text{league} & \text{Liverpool} \\ C & & & \\ D & & & \end{matrix}$$

- Find count matrix

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

- Normalize:

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

## 1.2 The Forward Algorithm

One type of inference problem that can be answered by an HMM is **Evaluation** - computing the probability of a sequence of observations. We calculate the likelihood of observing the validation sequence:

Liverpool win league

To do so, we calculate the forward probability matrix  $\alpha$ . Recall that

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

We have the following bottom-up dynamic programming algorithm to calculate the forward probabilities:

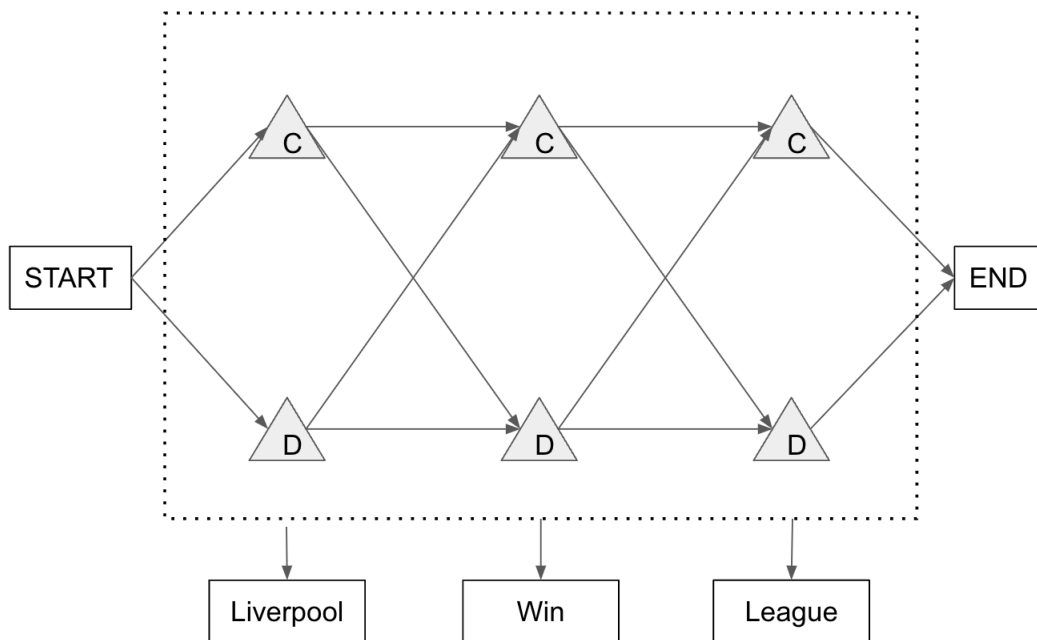
---

```

for t = 1, ..., T:
  for j = 1, ..., k:
    if t == 1:
       $\alpha_1(s_j) = \pi_j * A_{j,x_1}$ 
    else:
       $\alpha_t(s_j) = A_{j,x_t} * \alpha_{t-1}(s_k) * B_{k,j}$ 

```

---



$$\bullet \alpha_1 = \begin{bmatrix} \alpha_1(C) \\ \alpha_1(D) \end{bmatrix} = \begin{bmatrix} \pi_C * A_{C,x_1} \\ \pi_D * A_{D,x_1} \end{bmatrix} = \begin{bmatrix} \pi_C * A_{C,Liverpool} \\ \pi_D * A_{D,Liverpool} \end{bmatrix} = \begin{bmatrix} 3/5 \cdot 1/8 \\ 2/5 \cdot 2/3 \end{bmatrix} = \begin{bmatrix} 0.075 \\ 0.2667 \end{bmatrix}$$

- Using the path weight interpretation, we can write  $\alpha_2$  as the product of an edge weight matrix and  $\alpha_1$ . The edge weight  $w_{jkt}$  represents the weight on the edge in the trellis from hidden state  $k$  to hidden state  $j$  at time  $t$ , and can be calculated by

$$w_{jkt} = P(Y_t = s_j | Y_{t-1} = s_k) P(X_t = x_t | Y_t = s_j) = B_{kj} A_{jx_t}$$

$$\alpha_2 = \begin{bmatrix} w_{1,1,2} & w_{1,2,2} \\ w_{2,1,2} & w_{2,2,2} \end{bmatrix} \alpha_1 = \begin{bmatrix} 2/5 \cdot 1/2 & 3/4 \cdot 1/2 \\ 3/5 \cdot 1/6 & 1/4 \cdot 1/6 \end{bmatrix} \begin{bmatrix} 0.075 \\ 0.2667 \end{bmatrix} = \begin{bmatrix} 0.1150 \\ 0.0186 \end{bmatrix}$$

$$\bullet \alpha_3 = \begin{bmatrix} 0.0225 \\ 0.0123 \end{bmatrix} \alpha_2 = \begin{bmatrix} w_{1,1,3} & w_{1,2,3} \\ w_{2,1,3} & w_{2,2,3} \end{bmatrix} \alpha_2 = \begin{bmatrix} 2/5 \cdot 3/8 & 3/4 \cdot 3/8 \\ 3/5 \cdot 1/6 & 1/4 \cdot 1/6 \end{bmatrix} \begin{bmatrix} 0.1150 \\ 0.0186 \end{bmatrix} =$$

To find the likelihood of observing the validation sequence, all we need are the final forward probabilities:

$$\begin{aligned} & P(X_1 = \text{Liverpool}, X_2 = \text{win}, X_3 = \text{league}) \\ &= \sum_{y_3 \in \{C, D\}} P(x_1 = \text{Liverpool}, x_2 = \text{win}, x_3 = \text{league}, Y_3 = y_t) \\ &= \sum_{y_t \in \{C, D\}} \alpha_3(y_t) \\ &= 0.0225 + 0.0123 \\ &= 0.0348 \end{aligned}$$

### 1.3 The Backward Algorithm

Another type of inference problem that can be answered by an HMM is computing **Marginals** - computing the marginal probability distribution for a hidden state, given a sequence of observations. Recall that

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

Therefore, along with the forward probability matrix  $\alpha$ , we need to find the backward probability matrix  $\beta$ , where

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

We have a similar bottom-up dynamic programming algorithm to calculate the backward probabilities:

---

```

for t = T, ..., 1:
  for j = 1, ..., k:
    if t == T:
      beta_T(s_j) = 1
    else:
      beta_t(s_j) = A_{k,x_{t+1}} beta_{t+1}(s_k) B_{j,k}

```

---

- We can use the same path weight interpretation as for the forward algorithm, but for the backward algorithm the trellis edges are reversed. The edge weight  $w_{jkt}$  represents the weight on the edge in the trellis from hidden state  $k$  to hidden state  $j$  at time  $t$ , and can be calculated by

$$w_{jkt} = P(Y_{t+1} = s_k | Y_t = s_j) P(X_{t+1} = x_{t+1} | Y_{t+1} = s_k) = B_{jk} A_{kx_{t+1}}$$

$$\beta_2 = \begin{bmatrix} w_{1,1,2} & w_{1,2,2} \\ w_{2,1,2} & w_{2,2,2} \end{bmatrix} \beta_3 = \begin{bmatrix} 2/5 \cdot 3/8 & 3/5 \cdot 1/6 \\ 3/4 \cdot 3/8 & 1/4 \cdot 1/6 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.25 \\ 0.3229 \end{bmatrix}$$

$$\bullet \beta_1 = \begin{bmatrix} 0.0823 \\ 0.1072 \end{bmatrix} \beta_2 = \begin{bmatrix} w_{1,1,1} & w_{1,2,1} \\ w_{2,1,1} & w_{2,2,1} \end{bmatrix} \beta_2 = \begin{bmatrix} 2/5 \cdot 1/2 & 3/5 \cdot 1/6 \\ 3/4 \cdot 1/2 & 1/4 \cdot 1/6 \end{bmatrix} \begin{bmatrix} 0.25 \\ 0.3229 \end{bmatrix} =$$





Now, we have our  $\alpha$  and  $\beta$  matrices:

$$\alpha = \begin{array}{c} 1 \\ 2 \\ 3 \end{array} \begin{array}{cc} C & D \\ \left[ \begin{array}{cc} 0.0750 & 0.2667 \\ 0.1150 & 0.0186 \\ 0.0225 & 0.0123 \end{array} \right] \end{array}$$

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

1. What is  $P(Y_2 = C|\vec{x})$ ?

$$P(Y_2 = C|\vec{x}) = \frac{\alpha_2(C)\beta_2(C)}{P(\vec{x})} = \frac{0.1150 \times 0.2500}{0.0348} = 0.8261$$

2. What is  $P(Y_2 = D|\vec{x})$ ?

$$P(Y_2 = D|\vec{x}) = \frac{\alpha_2(D)\beta_2(D)}{P(\vec{x})} = \frac{0.0186 \times 0.3229}{0.0348} = 0.1726$$

3. What is  $P(Y_3 = C|\vec{x})$ ?

$$P(Y_2 = C|\vec{x}) = \frac{\alpha_2(C)\beta_2(C)}{P(\vec{x})} = \frac{0.0225 \times 1}{0.0348} = 0.6466$$

4. What is the minimum Bayes risk (MBR) decoder prediction for  $Y_2$ ?

This is the argmax over all marginal probabilities for the 2nd hidden state; here  $P(Y_2 = C|\vec{x}) > P(Y_2 = D|\vec{x})$  so the prediction is hidden state  $C$ .

## 1.4 The Viterbi Algorithm

Instead of finding the most likely hidden state at some time  $t$ , we may instead want to find the most likely sequence of hidden states. This is known as **Viterbi Decoding** - computing the most probable assignment of hidden states, given a sequence of observations.

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`

1. Recall that:

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

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

Using the formulae above and the first order Markov assumption, derive a recursive definition for  $\omega_t(s_k)$  and  $b_t(s_k)$  that will let you employ bottom-up dynamic programming.

First, we decompose the joint probability:

$$P(x_{1:t}, y_{1:t-1}, y_t) = P(x_t, x_{1:t-1}, y_t, y_{t-1}, y_{1:t-1})$$

$$= P(x_t|y_t)P(y_t|y_{t-1})P(x_{1:t-1}, y_{1:t-2}, y_{t-1})$$

(First order Markov Assumption)

Now, for maximizing:

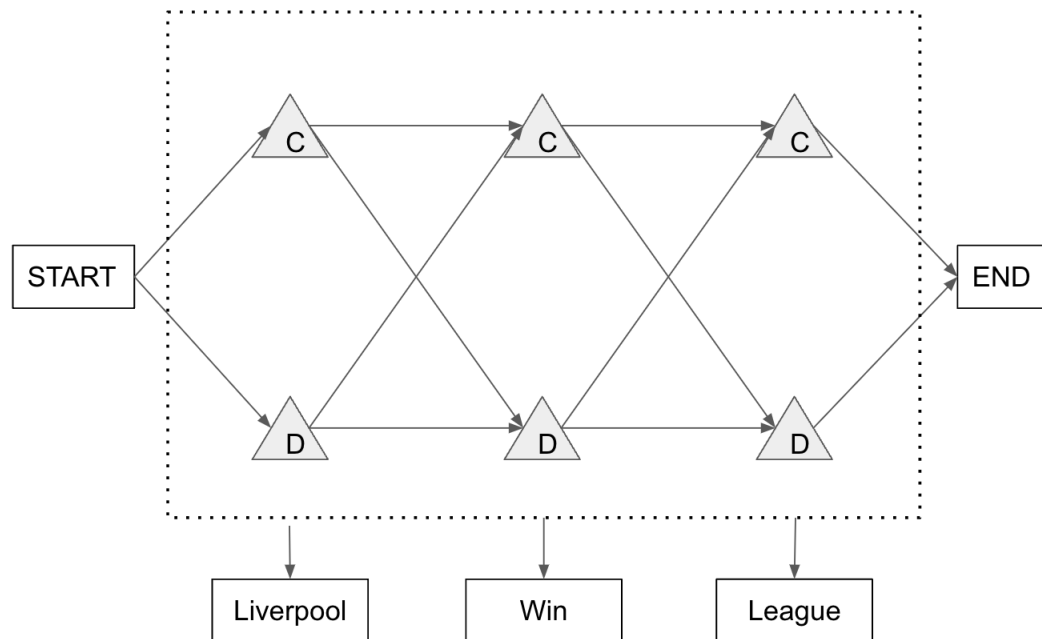
$$\omega_t(s_k) = \max_{s_j} (P(x_t|y_t = s_k)P(y_t = s_k|y_{t-1} = s_j) \max_{y_{1:t-2}} (P(x_{1:t-1}, y_{1:t-2}, y_{t-1} = s_j)))$$

$\implies$

$$\omega_t(s_k) = \max_{s_j} A_{s_k, x_t} \cdot B_{s_j, s_k} \cdot \omega_{t-1}(s_j)$$

$$b_t(s_k) = \arg \max_{s_j} A_{s_k, x_t} \cdot B_{s_j, s_k} \cdot \omega_{t-1}(s_j)$$

Below is the trellis corresponding to the given data:



2. Annotate the trellis at the nodes that correspond to:

- (a)  $\omega_1(C)$
- (b)  $\omega_1(D)$
- (c)  $\omega_2(C)$
- (d)  $\omega_2(D)$
- (e)  $\omega_3(C)$
- (f)  $\omega_3(D)$

3. Find the most likely sequence of tags given the observed data:

(a) Set up the matrices  $\omega$  and  $b$

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

and

$$b = \begin{array}{c} b_1 \\ b_2 \\ b_3 \\ b_4 \end{array} \begin{array}{c} \text{C} \quad \text{D} \quad \text{END} \\ \left[ \begin{array}{ccc} - & - & - \\ - & - & - \\ - & - & - \\ - & - & - \end{array} \right] \end{array}$$

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

(b) Solve for matrix entries using Dynamic Programming:

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

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

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

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

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

$$b_2(\mathbf{C}) = \mathbf{D}$$

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

$$b_2(\mathbf{D}) = \mathbf{D}$$

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

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

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

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

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

$$\hat{y}_{T+1} = \text{END}$$

$$\begin{aligned}\hat{y}_3 &= b_4(\text{END}) \\ &= \text{C}\end{aligned}$$

$$\begin{aligned}\hat{y}_2 &= b_3(\text{C} \quad ) \\ &= \text{C}\end{aligned}$$

$$\begin{aligned}\hat{y}_1 &= b_2(\text{C} \quad ) \\ &= \text{D}\end{aligned}$$

$$\begin{aligned}\hat{y}_0 &= b_1(\text{D} \quad ) \\ &= \text{START}\end{aligned}$$

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

## 2 Working in Log-space

### 2.1 Motivation

Some of the probabilities we work with in Homework 7 about are tiny and some of them are much larger. We tend to work with the tiny ones in log-space and only get back probabilities if we really need them for some other purpose. Throughout HW7 you will keep your probabilities in log-space. In this section we will motivate why we use log-space for small values.

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?

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

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

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?

$$\log(P(x_1 = 1, x_2 = 1, x_3 = 1)) \\ = \log(x_1 = 1) + \log(P(x_2 = 1 \mid x_1 = 1)) + \log(P(x_3 = 1 \mid x_2 = 1, x_1 = 1)) \\ = \log(0.002) + \log(0.004) + \log(0.003) \\ = -6.2146 - 5.8091 - 5.5215 \\ = -17.5452$$

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 entries in  $\alpha$  can be computed using the bottom-up dynamic programming algorithm:

- $\alpha_1(j) = \pi_j A_{jx_1}$
  - 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  $\log A_{kj}$ 

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

In the backward algorithm, we also have a similar bottom-up dynamic programming algorithm:

- $\beta_T(j) = 1$
  - 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  $\log(B_{jk})$ 

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

After transforming the equations into log form, you may discover calculations of the following type:

$$\log \sum_i \exp(v_i)$$

This may be programmed as is, but  $\exp(v_i)$  may cause underflow when  $v_i$  is large and negative. One way to avoid this is to use the [log-sum-exp trick](#).

The log-sum-exp trick simply adds the maximum value in the vector to the log probabilities as follows:

$$\log \sum_i \exp(v_i - m) + \max_i(v_i)$$

## 3 Dynamic Programming

### [DP Notebook](#)

To access this Colab notebook you will need to be logged into Google Drive with your Andrew email.