

# COMS W4701 HW4

Dawei He (dh3027)

April 10, 2022

## Problem 1

(a)

By conditioning, we have

$$Pr(X_t|x_k, x_j) = \frac{Pr(X_t, x_k, x_j)}{Pr(x_k, x_j)}$$

because the denominator  $Pr(x_k, x_j)$  is constant, so we can say  $Pr(X_t|x_k, x_j)$  is proportional to  $Pr(X_t, x_k, x_j)$ , which can be represent as

$$Pr(X_t|x_k, x_j) \propto Pr(X_t, x_k, x_j)$$

(b)

We need to sum over all  $x$  except  $x_t, x_k, x_j$ , i.e.,

$$Pr(x_t, x_k, x_j) = \sum_{x_i, i=0, \dots, t \text{ \& } i \neq j, k, t} Pr(x_0, x_1, \dots, x_t)$$

.

(c)

Without losing generalization,

$$\begin{aligned} & Pr(x_t, x_k, x_j) \\ &= \sum_{x_i, i=0, \dots, t \text{ \& } i \neq j, k, t} Pr(x_0, x_1, \dots, x_t) \\ &= \sum_{x_i, i=0, \dots, t \text{ \& } i \neq j, k, t} Pr(x_0) \prod_{i=1}^t Pr(x_i|x_{i-1}) \\ &= \sum_{x_i, i=k+1, \dots, t-1} \sum_{x_i, i=j+1, \dots, k-1} \sum_{x_i, i=0, \dots, j-1} Pr(x_0) \prod_{i=1}^t Pr(x_i|x_{i-1}) \end{aligned}$$

.

Consider the product  $Pr(x_0) \prod_{i=1}^j Pr(x_i|x_{i-1})$ , which is equal to

$$Pr(x_0)Pr(x_1|x_0)Pr(x_2|x_1) \cdots Pr(x_j|x_{j-1}) = Pr(x_0, x_1)Pr(x_2|x_1) \cdots Pr(x_j|x_{j-1}),$$

since  $Pr(x_i|x_{i-1}) = Pr(x_i|x_{i-1}, \dots, x_0)$  in Markov Chain, then  $Pr(x_2|x_1) = Pr(x_2|x_1, x_0)$  etc. The product can be written as

$$Pr(x_0) \prod_{i=1}^j Pr(x_i|x_{i-1}) = Pr(x_0, x_1, \dots, x_j).$$

And  $\sum_{x_i, i=0, \dots, j-1}$  means summing the first  $j$  elements over  $Pr(x_0, x_1, \dots, x_j)$ , then get the marginal distribution of  $x_j$ , i.e.,

$$\sum_{x_i, i=0, \dots, j-1} Pr(x_0) \prod_{i=1}^j Pr(x_i|x_{i-1}) = Pr(x_j).$$

Now consider the product

$$\prod_{i=j+1}^k Pr(x_i|x_{i-1}) = Pr(x_{j+1}|x_j)Pr(x_{j+2}|x_{j+1}) \cdots Pr(x_k|x_{k-1}),$$

since

$$\begin{aligned} & Pr(x_{j+1}|x_j)Pr(x_{j+2}|x_{j+1}) \\ &= \frac{Pr(x_j)Pr(x_{j+1}|x_j)Pr(x_{j+2}|x_{j+1})}{Pr(x_j)} \\ &= \frac{Pr(x_j)Pr(x_{j+1}|x_j)Pr(x_{j+2}|x_{j+1}, x_j)}{Pr(x_j)} \\ &= Pr(x_{j+1}, x_{j+2}|x_j), \end{aligned}$$

the product can be reduced as  $\prod_{i=j+1}^k Pr(x_i|x_{i-1}) = Pr(x_{j+1}, x_{j+2}, \dots, x_k|x_j)$ ,  $\sum_{x_i, i=j+1, \dots, k-1}$  means summing the middle  $k-j-1$  elements over  $Pr(x_{j+1}, x_{j+2}, \dots, x_k|x_j)$ , then get the marginal distribution given  $x_j$ , i.e.,

$$\sum_{x_i, i=j+1, \dots, k-1} \prod_{i=j+1}^k Pr(x_i|x_{i-1}) = \sum_{x_i, i=j+1, \dots, k-1} Pr(x_{j+1}, x_{j+2}, \dots, x_k|x_j) = Pr(x_k|x_j).$$

Then consider the product

$$\prod_{i=k+1}^t Pr(x_i|x_{i-1}) = Pr(x_{k+1}|x_k)Pr(x_{k+2}|x_{k+1}) \cdots Pr(x_t|x_{t-1}).$$

same as before, the product can be reduced as marginal distribution given  $x_t$ ,  $Pr(x_t|x_k)$ , when summing over  $x_i$  from  $i = k+1$  to  $i = t-1$ .

Hence the formula the products of the three parts:

$$\begin{aligned} & Pr(x_t, x_k, x_j) \\ &= \sum_{x_i, i=k+1, \dots, t-1} \sum_{x_i, i=j+1, \dots, k-1} \sum_{x_i, i=0, \dots, j-1} Pr(x_0) \prod_{i=1}^t Pr(x_i|x_{i-1}) \\ &= Pr(x_j)Pr(x_k|x_j)Pr(x_t|x_k). \end{aligned}$$

(d)

According to part(c),

$$Pr(x_t|x_k, x_j) = \frac{Pr(x_t, x_k, x_j)}{Pr(x_k, x_j)} = \frac{Pr(x_j)Pr(x_k|x_j)Pr(x_t|x_k)}{Pr(x_k, x_j)} = \frac{Pr(x_k, x_j)Pr(x_t|x_k)}{Pr(x_k, x_j)} = Pr(x_t|x_k),$$

and formula in part(c) holds true for all values of  $x_t$ , which means  $Pr(X_t|x_k, x_j) = Pr(X_t|x_k)$  is true. Hence,  $X_t$  is conditionally independent of  $X_j$  and  $X_k$ .

## Problem 2

(a)

The stochastic matrix  $T$  is

$$T = \begin{bmatrix} 0 & 1/3 & 0 & 1/2 & 0 & 0 \\ 1/2 & 0 & 1 & 0 & 1/3 & 0 \\ 0 & 1/3 & 0 & 0 & 0 & 0 \\ 1/2 & 0 & 0 & 0 & 1/3 & 0 \\ 0 & 1/3 & 0 & 1/2 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1/3 & 0 \end{bmatrix}$$

We have to compute the eigenvector with corresponding eigenvalue being 1 and normalize this eigenvector so all elements sum up to 1. Using numpy to compute above procedures and we get the stationary distribution  $\pi$  over the robot's predicted location:

$$\pi = [P(A) \ P(B) \ P(C) \ P(D) \ P(E) \ P(F)] = [1/6 \ 1/4 \ 1/12 \ 1/6 \ 1/4 \ 1/12]$$

(b)

The beginning state is  $P(X_0) = (0, 0, 0, 0, 1, 0)$ , we can get the probability of state  $X_1$  by calculating

$$P(X_1) = T \cdot P(X_0) = (0, 1/3, 0, 1/3, 0, 1/3)$$

now we have to compute the probability of state  $X_1$  given  $e_1$ .

$$P(X_1|e_1) = \frac{P(e_1|X_1)P(X_1)}{P(e_1)} = \frac{1}{1/2 \times 0 + 1/4 \times 1/3 + 1/2 \times 0 + 1/4 \times 1/3 + 0 \times 0 + 1/4 \times 1/3} \begin{bmatrix} 1/2 \times 0 \\ 1/4 \times 1/3 \\ 1/2 \times 0 \\ 1/4 \times 1/3 \\ 0 \times 0 \\ 1/4 \times 1/3 \end{bmatrix} = \begin{bmatrix} 0 \\ 1/3 \\ 0 \\ 1/3 \\ 0 \\ 1/3 \end{bmatrix}$$

(c)

$$P(X_2|e_1) = \alpha'_2 = T \begin{bmatrix} 0 \\ 1/3 \\ 0 \\ 1/3 \\ 0 \\ 1/3 \end{bmatrix} = \begin{bmatrix} 5/18 \\ 0 \\ 1/9 \\ 0 \\ 11/18 \\ 0 \end{bmatrix}$$

$$\alpha_2 = O_2 \begin{bmatrix} 5/18 \\ 0 \\ 1/9 \\ 0 \\ 11/18 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \times 5/18 \\ 0 \times 0 \\ 1/4 \times 1/9 \\ 1/4 \times 0 \\ 1/4 \times 11/18 \\ 1/2 \times 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1/36 \\ 0 \\ 11/72 \\ 0 \end{bmatrix}$$

$$P(X_2|e_1, e_2) = \frac{P(e_2|X_2, e_1)P(X_2|e_1)}{P(e_2|e_1)} = \frac{1}{1/36 + 11/72} \begin{bmatrix} 0 \\ 0 \\ 1/36 \\ 0 \\ 11/72 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 2/13 \\ 0 \\ 11/13 \\ 0 \end{bmatrix}$$

(d)

According to (b),

$$P(X_1, e_1) = P(e_1|X_1)P(X_1) = \begin{bmatrix} 0 \\ 1/12 \\ 0 \\ 1/12 \\ 0 \\ 1/12 \end{bmatrix}$$

we can get calculate  $P(X_1, X_2, e_1, e_2)$

$$\begin{aligned} P(X_1, X_2, e_1, e_2) &= P(e_2|X_2)P(X_2|X_1)P(X_1, e_1) \\ &= P(e_2|X_2) \begin{bmatrix} 0 & 1/36 & 0 & 1/24 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1/36 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1/36 & 0 & 1/24 & 0 & 1/12 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1/144 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1/144 & 0 & 1/96 & 0 & 1/48 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{aligned}$$

finally we can get  $P(X_1, X_2|e_1, e_2)$

$$P(X_1, X_2|e_1, e_2) = \frac{P(X_1, X_2, e_1, e_2)}{P(e_1, e_2)} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2/13 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2/13 & 0 & 3/13 & 0 & 6/13 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

The most likely state sequence is  $E \rightarrow F \rightarrow E$

Note:  $P(X_2|X_1)P(X_1, e_1)$  is row element wise multiplication,  $P(e_2|X_2) \cdot [P(X_2|X_1)P(X_1, e_1)]$  is column element wise multiplication.

(e)

$P(e_2|X_1)$  is  $\beta_1$ , we want to use backward algorithm to compute it. First we initialize

$$\beta_2 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

Then we can compute  $\beta'_2$

$$\beta'_2 = O_2\beta_3 = \begin{bmatrix} 0 \times 1 \\ 0 \times 1 \\ 1/4 \times 1 \\ 1/4 \times 1 \\ 1/4 \times 1 \\ 1/2 \times 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1/4 \\ 1/4 \\ 1/4 \\ 1/2 \end{bmatrix}$$

We can get  $\beta_1$

$$P(e_2|X_1) = \beta_1 = T^T \beta'_2 = \begin{bmatrix} 1/8 \\ 1/6 \\ 0 \\ 1/8 \\ 1/4 \\ 1/4 \end{bmatrix}$$

The quantity means if we only know we are at state  $X_1$  and don't know the observation  $e_1$  at state  $X_1$ , what is the probability that we can observe # at state  $X_2$ .

(f)

$$P(X_1, e_2|e_1) = P(X_1|e_1)P(e_2|X_1) = \begin{bmatrix} 1/8 \\ 1/6 \\ 0 \\ 1/8 \\ 1/4 \\ 1/4 \end{bmatrix} \times \begin{bmatrix} 0 \\ 1/3 \\ 0 \\ 1/3 \\ 0 \\ 1/3 \end{bmatrix} = \begin{bmatrix} 0 \\ 1/18 \\ 0 \\ 1/24 \\ 0 \\ 1/12 \end{bmatrix}$$

$$P(X_1|e_2, e_1) = \frac{1}{1/18 + 1/24 + 1/12} \begin{bmatrix} 0 \\ 1/18 \\ 0 \\ 1/24 \\ 0 \\ 1/12 \end{bmatrix} = \begin{bmatrix} 0 \\ 4/13 \\ 0 \\ 3/13 \\ 0 \\ 6/13 \end{bmatrix}$$

According to (d)

$$P(X_1, X_2|e_1, e_2) = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2/13 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2/13 & 0 & 3/13 & 0 & 6/13 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

So we can marginalize  $X_2$

$$P(X_1|e_2, e_1) = \sum_{X_2} P(X_1, X_2|e_1, e_2) = [0 \quad 4/13 \quad 0 \quad 3/13 \quad 0 \quad 6/13]$$

We can see two results of two methods are the same.

## Problem 3

### 3.3

(a)

Train accuracy: 0.9555019712104628

Test accuracy: 0.8375544528500484

Because some words in testing dataset do not occur in the training dataset, the observation probability  $P(e|x)$  will be given 1 for all POS state. So their state will be less accurate than those words observed in the training set. So the overall test accuracy will be less than training accuracy

(b)

There could be some small portion of invalid sentences in the training dataset, which means these sentence could have wrong sequence of POS. HMM will maximize the likelihood of POS for most of the correct sentences in the training set, and output their correct sequence of states. But if the sentence has wrong POS sequence, HMM can not output such wrong states by back propagation.

### 3.5

(c)

```
{'What': array([0.          , 0.          , 0.          , 0.          , 0.          ,
0.00811993, 0.          , 0.          , 0.          , 0.          ,
0.9282028 , 0.          , 0.          , 0.          , 0.          ,
0.          , 0.          ]), 'GoogleOS': array([0.71930156, 0.01736604, 0.97535508, 0.99477996, 0.87332341,
0.02439862, 0.93150538, 0.95648066, 0.63177494, 0.82638099,
0.06247148, 0.05776416, 0.04828422, 0.00488973, 0.52776901,
0.9901241 , 0.95996977]), 'Google': array([0.          , 0.          , 0.          , 0.          , 0.          ,
0.          , 0.          , 0.          , 0.          , 0.          ,
0.          , 0.64093566, 0.          , 0.          , 0.          ,
0.          , 0.          ]), 'Morphed': array([0.28069844, 0.04616681, 0.02464492, 0.00522004, 0.12667659,
0.96748145, 0.06849462, 0.04351934, 0.36822506, 0.17361901,
0.00932571, 0.30130018, 0.02702465, 0.00187195, 0.47223099,
0.0098759 , 0.04003023]), 'Into': array([0.          , 0.92985294, 0.          , 0.          , 0.          ,
0.          , 0.          , 0.          , 0.          , 0.          ,
0.          , 0.          , 0.          , 0.          , 0.          ,
0.          , 0.          ]), 'if': array([0.          , 0.0066142 , 0.          , 0.          , 0.          ,
0.          , 0.          , 0.          , 0.          , 0.          ,
0.          , 0.          , 0.99323832, 0.          , 0.          ,
0.          , 0.          ]), '?': array([0.          , 0.          , 0.          , 0.          , 0.          ,
0.          , 0.          , 0.          , 0.          , 0.          ,
0.          , 0.          , 0.92469113, 0.          , 0.          ,
0.          , 0.          ])}
```

Figure 1: Observation probabilities for 3.5

$\arg \max_X P(e = \textit{What}|X) = \text{PRON}$ , which is correct.  
 $\arg \max_X P(e = \textit{GoogleOS}|X) = \text{AUX}$ , which is false.  
 $\arg \max_X P(e = \textit{Google}|X) = \text{PROPN}$ , which is correct.  
 $\arg \max_X P(e = \textit{Morphed}|X) = \text{CCONJ}$ , which is false.  
 $\arg \max_X P(e = \textit{Into}|X) = \text{ADP}$ , which is correct.  
 $\arg \max_X P(e = \textit{if}|X) = \text{SCONJ}$ , which is correct.  
 $\arg \max_X P(e = \textit{?}|X) = \text{PUNCT}$ , which is correct.

POS of *GoogleOS* and *Morphed* are hard to identify and appear to be ambiguous. They both do not appear in the training dataset and *GoogleOS* is a compound word, so it is hard for model to correctly produce their POS.