

# HOMWORK 7: HIDDEN MARKOV MODELS

10-301/10-601 Introduction to Machine Learning (Fall 2022)

<http://www.cs.cmu.edu/~mgormley/courses/10601/>

OUT: Friday, November 11th

DUE: Monday, November 21st

TAs: Monica, Shubham, Kalvin, Brandon, Chu

**Summary** In this assignment you will implement a new named entity recognition system using Hidden Markov Models. You will begin by going through some multiple choice and short answer warm-up problems to build your intuition for these models and then use that intuition to build your own HMM models.

## START HERE: Instructions

- **Collaboration Policy:** Please read the collaboration policy here: <http://www.cs.cmu.edu/~mgormley/courses/10601/syllabus.html>
- **Late Submission Policy:** See the late submission policy here: <http://www.cs.cmu.edu/~mgormley/courses/10601/syllabus.html>
- **Submitting your work:** You will use Gradescope to submit answers to all questions and code. Please follow instructions at the end of this PDF to correctly submit all your code to Gradescope.
  - **Written:** For written problems such as short answer, multiple choice, derivations, proofs, or plots, please use the provided template. Submissions can be handwritten onto the template, but should be labeled and clearly legible. If your writing is not legible, you will not be awarded marks. Alternatively, submissions can be written in  $\LaTeX$ . Each derivation/proof should be completed in the boxes provided. You are responsible for ensuring that your submission contains exactly the same number of pages and the same alignment as our PDF template. If you do not follow the template, your assignment may not be graded correctly by our AI assisted grader and there will be a **5% penalty** (e.g., if the homework is out of 100 points, 5 points will be deducted from your final score).
  - **Programming:** You will submit your code for programming questions on the homework to Gradescope (<https://gradescope.com>). After uploading your code, our grading scripts will autograde your assignment by running your program on a virtual machine (VM). When you are developing, check that the version number of the programming language environment (e.g. Python 3.9.12) and versions of permitted libraries (e.g. `numpy` 1.23.0) match those used on Gradescope. You have 10 free Gradescope programming submissions. After 10 submissions, you will begin to lose points from your total programming score. We recommend debugging your implementation on your local machine (or the Linux servers) and making sure your code is running correctly first before submitting your code to Gradescope.
- **Materials:** The data and reference output that you will need in order to complete this assignment is posted along with the writeup and template on the course website.

## Instructions for Specific Problem Types

For “Select One” questions, please fill in the appropriate bubble completely:

**Select One:** Who taught this course?

- ☒ Matt Gormley
- ☐ Marie Curie
- ☐ Noam Chomsky

If you need to change your answer, you may cross out the previous answer and bubble in the new answer:

**Select One:** Who taught this course?

- ☒ Henry Chai
- ☐ Marie Curie
- ☒ Noam Chomsky

For “Select all that apply” questions, please fill in all appropriate squares completely:

**Select all that apply:** Which are scientists?

- ☒ Stephen Hawking
- ☒ Albert Einstein
- ☒ Isaac Newton
- ☐ I don't know

Again, if you need to change your answer, you may cross out the previous answer(s) and bubble in the new answer(s):

**Select all that apply:** Which are scientists?

- ☒ Stephen Hawking
- ☒ Albert Einstein
- ☒ Isaac Newton
- ☒ I don't know

For questions where you must fill in a blank, please make sure your final answer is fully included in the given space. You may cross out answers or parts of answers, but the final answer must still be within the given space.

**Fill in the blank:** What is the course number?

10-601

10-~~6~~301

## Written Questions (49 points)

### 1 L<sup>A</sup>T<sub>E</sub>X Bonus Point (1 points)

1. (1 point) **Select one:** Did you use L<sup>A</sup>T<sub>E</sub>X for the entire written portion of this homework?

☒ Yes

☐ No

### 2 Hidden Markov Models (12 points)

1. In this section, we will test your understanding of several aspects of HMMs.

- (a) (2 points) **Select all that apply:** Let  $Y_t$  be the state at time  $t$ . Which of the following are true under the (first-order) Markov assumption in an HMM?

☐ The states are independent

☐ The observations are independent

☐  $Y_t \perp\!\!\!\perp Y_{t-1} \mid Y_{t-2}$

☒  $Y_t \perp\!\!\!\perp Y_{t-2} \mid Y_{t-1}$

☐ None of the above

- (b) (2 points) **Select all that apply:** Which of the following independence assumptions hold in an HMM?

☒ The current observation  $X_t$  is conditionally independent of all other observations given the current state  $Y_t$

☒ The current observation  $X_t$  is conditionally independent of all other states given the current state  $Y_t$

☐ The current state  $Y_t$  is conditionally independent of all states given the previous state

☐ The current observation  $X_t$  is conditionally independent of  $Y_{t-2}$  given the previous observation  $X_{t-1}$

☐ None of the above

2. In the remaining questions, you will see two quantities and decide what the strongest relation between them is. For each question, there is **only one correct answer**. As a reminder,  $\alpha_t(s_j) = P(Y_t = s_j, x_{1:t})$  and  $\beta_t(s_j) = P(x_{t+1:T} \mid Y_t = s_j)$ . Assume that there are  $J$  possible hidden states, so each  $Y_t \in \{s_i\}_{i=1}^J$ , and that we do not use pseudocounts.

*Note:* ? means it's not possible to assign any true relation.

- (a) (3 points) **Select one:** What is the relation between  $\sum_{i=1}^J (\alpha_5(s_i)\beta_5(s_i))$  and  $P(x_{1:T})$ , for  $T > 5$ ? Select only the **strongest** relation that necessarily holds.

☒ =

☐ >

☐ <

☐  $\leq$

☐  $\geq$

☐ ?

- (b) (3 points) **Select one:** What is the relation between  $P(Y_4 = s_1, Y_5 = s_2, x_{1:T})$  and  $\alpha_4(s_1)\beta_5(s_2)$ ? Select only the **strongest** relation that necessarily holds.

☐ =

☐ >

☐ <

☒  $\leq$

☐  $\geq$

☐ ?

- (c) (2 points) **Select one:** What is the relation between  $\alpha_5(s_i)$  and  $\beta_5(s_i)$ ? Select only the **strongest** relation that necessarily holds.

☐ =

☐ >

☐ <

☐  $\leq$

☐  $\geq$

☒ ?

### 3 Viterbi Decoding (4 points)

Suppose we have a set of sequences consisting of  $T$  observed states,  $x_1, \dots, x_T$ , where each  $x_t \in \{1, 2, 3\}$ . Each observed state is associated with a hidden state  $Y_t \in \{C, D\}$ . Let  $s_1 = C$  and  $s_2 = D$ .

In the Viterbi algorithm, we seek to find the most probable hidden state sequence  $\hat{y}_1, \dots, \hat{y}_T$  given the observations  $x_1, \dots, x_T$ .

We define:

- $\mathbf{B}$  to be the transition matrix:  $B_{jk} = P(Y_t = s_k \mid Y_{t-1} = s_j)$
- $\mathbf{A}$  to be the emission matrix:  $A_{jk} = P(X_t = k \mid Y_t = s_j)$
- $\pi$  to be the initialization matrix for  $Y_1$ :  $\pi_j = P(Y_1 = s_j)$
- $\omega_t(s_k)$  to be the maximum product of all the probabilities taken through path  $Y_1, \dots, Y_{t-1}$  that ends with  $Y_t$  at state  $s_k$ .

$$\omega_t(s_k) = \max_{y_1, \dots, y_{t-1}} P(x_{1:t}, y_{1:t-1}, Y_t = s_k) \quad (1)$$

- $b_t(s_k)$  to be the backpointer that stores the path through hidden states that gives us the highest product.

$$b_t(s_k) = \operatorname{argmax}_{y_1, \dots, y_{t-1}} P(x_{1:t}, y_{1:t-1}, Y_t = s_k) \quad (2)$$

We outline the Viterbi Algorithm below:

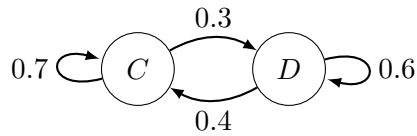
1. Initialize  $\omega_1(s_j) = \pi_j A_{jx_1}$  and  $b_1(s_j) = j$
2. For  $t > 1$ , we have

$$\begin{aligned} \omega_t(s_j) &= \max_{k \in \{1, \dots, J\}} A_{jx_t} B_{kj} \omega_{t-1}(s_k) \\ b_t(s_j) &= \operatorname{argmax}_{k \in \{1, \dots, J\}} A_{jx_t} B_{kj} \omega_{t-1}(s_k) \end{aligned}$$

We can obtain the most probable sequence by backtracing through the backpointers as follows:

1.  $\hat{y}_T = \operatorname{argmax}_{k \in \{1, \dots, J\}} \omega_T(s_k)$ .
2. For  $t = T - 1, \dots, 1$ :  
 $\hat{y}_t = b_{t+1}(\hat{y}_{t+1})$
3. Return  $\hat{y}_1, \dots, \hat{y}_T$

For the following question, consider the Hidden Markov Model specified below. We have that  $P(Y_1 = C) = P(Y_1 = D) = 0.5$ , and the state transition model and emission probability tables are given as follows.



$k$	$P(X_t = k \mid Y_t = C)$	$P(X_t = k \mid Y_t = D)$
1	0.5	0.3
2	0.4	0.5
3	0.1	0.2

We observed  $X_1 = 1$  and  $X_2 = 2$ .

1. (2 points) Compute  $\omega_1(C)$  and  $\omega_1(D)$ . If your answers involve decimal numbers, please round your answer to **TWO** decimal places.

*Note:* Showing your work in these questions is optional, but it is recommended to help us understand where any misconceptions may occur. Only your answer in the left box will be graded.

What is  $\omega_1(C)$ ?

$\omega_1(C)$	Work
0.25	

What is  $\omega_1(D)$ ?

$\omega_1(D)$	Work
0.15	

2. (2 points) (**Select one**) Which of the following is the most likely sequence of hidden states?

- ☒  $Y_1 = C, Y_2 = C$
- ☐  $Y_1 = D, Y_2 = D$
- ☐  $Y_1 = D, Y_2 = C$
- ☐  $Y_1 = C, Y_2 = D$
- ☐ Not enough information.

## 4 Forward-Backward Algorithm (26 points)

The following questions should be completed before you start the programming component of this assignment. To help you prepare to implement the HMM forward-backward algorithm (see Section 7.5 for a detailed explanation), we have provided a small example for you to work through by hand. This toy data set consists of a training set of three sequences with three unique words and two tags, and a validation set with a single sequence composed of the same words occurring in the training set.

### Training set:

you        D  
eat        C  
fish       D

you        D  
fish       D  
eat        C

eat        C  
fish       D

The training word sequences are:

$$\mathbf{x}^{(1)} = [\text{you} \text{ eat} \text{ fish}]^T$$

$$\mathbf{x}^{(2)} = [\text{you} \text{ fish} \text{ eat}]^T$$

$$\mathbf{x}^{(3)} = [\text{eat} \text{ fish}]^T$$

And the corresponding tags are:

$$\mathbf{y}^{(1)} = [D \ C \ D]^T$$

$$\mathbf{y}^{(2)} = [D \ D \ C]^T$$

$$\mathbf{y}^{(3)} = [C \ D]^T$$

### Validation set:

fish  
eat  
you

The validation word sequence is:

$$\mathbf{x}_{\text{validation}} = [\text{fish} \text{ eat} \text{ you}]^T$$

In this question, we define:

- Each observed state  $x_t \in \{1, 2, 3\}$ , where 1 corresponds to `you`, 2 corresponds to `eat`, and 3 corresponds to `fish`
- Each hidden state  $Y_t \in \{C, D\}$ . Let  $s_1 = C$  and  $s_2 = D$ .
- $\mathbf{B}$  is the transition matrix, where  $B_{jk} = P(Y_t = s_k \mid Y_{t-1} = s_j)$ . Here  $\mathbf{B}$  is a  $2 \times 2$  matrix.
- $\mathbf{A}$  is the emission matrix, where  $A_{jk} = P(X_t = k \mid Y_t = s_j)$ . Here  $\mathbf{A}$  is a  $2 \times 3$  matrix. As an example,  $A_{23}$  denotes  $P(X_t = 3 \mid Y_t = s_2)$ , or the probability that  $X_t$  corresponds to `fish` given the hidden state  $Y_t = D$ .
- $\pi$  describes  $Y_1$ 's initialization probabilities:  $\pi_j = P(Y_1 = s_j)$ .

For pseudo-code of the Forward-Backward Algorithm, refer to [7.5](#).

*Note:* Pseudocounts used in section [7.4](#) should also be used here.

For all numerical answers, round to **four decimal places** after the decimal point. Showing your work in these questions is optional, but it is recommended to help us understand where any misconceptions may occur. Only your answer in the left box will be graded.



1. (5 points) Compute  $\alpha_2(C)$ , the  $\alpha$  value associated with the tag “C” for the second word in the validation sequence.

$\alpha_2(C)$	Work
0.1311	

2. (3 points) Compute  $\beta_2(D)$ , the  $\beta$  value associated with the tag “D” for the second word in the validation sequence.

$\beta_2(D)$	Work
0.2500	

3. (3 points) Predict the tag for the third word in the validation sequence.

Tag	Work
D	

4. (3 points) Compute the log-likelihood for the entire validation sequence, “fish eat you”, using  $e$  as the log base.

Log-Likelihood	Work
-3.0439	

5. (6 points) In the forward algorithm, recall that the entries in  $\alpha$  can be computed using the following dynamic programming algorithm:

```

 $\alpha_1(s_j) = \pi_j \cdot A_{s_j, x_1}$  for all  $s_j$ 
For  $t = 1, 2, 3, \dots, T$ 
     $\alpha_t(s_j) = A_{s_j, x_t} \sum_k B_{s_k, s_j} \alpha_{t-1}(s_k)$  for all  $s_j$ 

```

To implement the forward algorithm in log-space, we can derive  $\log(\alpha_t(s_j))$  in terms of  $\log(\alpha_{t-1})$ ,  $\log \mathbf{B}$  and  $\log \mathbf{A}$ :

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

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

```

 $\beta_T(s_j) = 1$  for all  $s_j$ 
For  $t = T - 1, T - 2, T - 3, \dots, 1$ 
     $\beta_t(s_j) = \sum_k \beta_{t+1}(s_k) A_{s_k, x_{t+1}} B_{s_j, s_k}$ 

```

We can also derive the backward algorithm in log-space by finding  $\log(\beta_t(j))$  in terms of  $\log(\beta_{t+1})$ ,  $\log \mathbf{B}$  and  $\log \mathbf{A}$ :

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

Note how the above equations include terms that look like  $\log \sum_i \exp(v_i)$ . In your programming section, you shouldn't program this as-is, because the  $\exp$  term will underflow when  $v_i$  is very negative. You'll instead implement this using a trick known as the log-sum-exp trick. For this written problem, you may assume access to a `logsumexp` function that takes in a matrix  $H$  and returns a vector  $v$  where  $v_i = \log \sum_j \exp(H_{ij})$ . In other words, `logsumexp` performs the log-sum-exp trick on each row of  $H$ .

Now it's time to vectorize the log-space equations above! For each of the following quantities, write one line of code to compute them from the given quantities. You are given that  $\log \alpha_t$  is a column vector which is the element-wise log of  $\alpha_t$ ,  $\log \beta_t$  is a column vector which is the element-wise log of

$\beta_t$ ,  $\log A$  is the element-wise log of the emission matrix,  $\log B$  is the element-wise log of the transition matrix,  $\log \pi$  is the element-wise log of the initial probabilities,  $\mathbf{x}$  is the input sequence (a list of words represented as matrix indices), and the `logsumexp` function as described above. You may only write your answer in terms of these quantities. You may assume that NumPy has been imported as `np` and that broadcasting is allowed (e.g. if  $M$  is an  $a \times b$  matrix and  $v$  is a vector of length  $b$ , then  $M + v$  will add  $v$  element-wise to every row of  $M$ ).

(a) (2 points)  $\log \alpha_1$

Your Answer

```
np.log( $\pi$ ) + np.logA[:, $\mathbf{x}[0]$ ]
```

(b) (2 points)  $\log \alpha_2$

Your Answer

```
np.logA[:, $\mathbf{x}[1]$ ] + logsumexp(np.log $\alpha_1$  + np.log $B.T$ )
```

(c) (2 points)  $\log \beta_2$

Your Answer

```
logsumexp(np.logA[:, $\mathbf{x}[2]$ ] + np.log $\beta_3$  + np.log $B$ )
```

## 5 Empirical Questions (6 points)

Return to these questions after implementing your `learnhmm.py` and `forwardbackward.py` functions. Please ensure that you have used the log-sum-exp trick in your programming as described in Section 7.5.3 before answering these empirical questions.

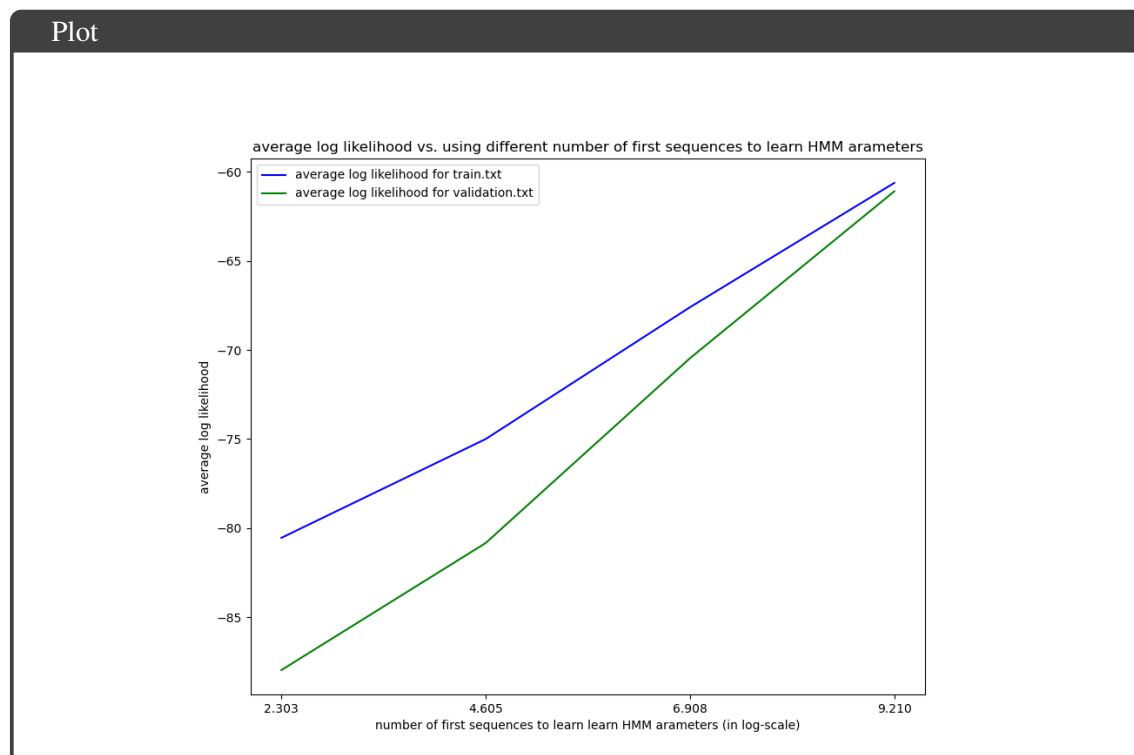
Using the full data set **en\_data/train.txt** in the handout, use your implementation of `learnhmm.py` to learn HMM parameters using the first 10, 100, 1000, and 10000 sequences in the file. Use these learned parameters to perform prediction on both the English **train.txt** and the **validation.txt** files with your implementation of `forwardbackward.py`. Construct a plot with number of sequences used for training on the x-axis and average log likelihood across all sequences from the English **train.txt** and the **validation.txt** on the y-axis (see Section 7.5 for details on computing the log data likelihood for a sequence). Please use **log-scale** with base  $e$  for the x-axis.

Fill in the table with the resulting log likelihood values, rounded to two decimal places, and include your plot in the large box. To receive credit for your plot, you must submit a computer generated plot. **DO NOT** hand draw your plot.

1. (4 points) Fill in this table.

# Sequences	Train Average Log-Likelihood	Validation Average Log-Likelihood
10	-80.54	-87.97
100	-74.99	-80.83
1000	-67.59	-70.46
10000	-60.61	-61.09

2. (2 points) Put your plot below:



## 6 Collaboration Questions

After you have completed all other components of this assignment, report your answers to these questions regarding the collaboration policy. Details of the policy can be found [here](#).

1. Did you receive any help whatsoever from anyone in solving this assignment? If so, include full details.
2. Did you give any help whatsoever to anyone in solving this assignment? If so, include full details.
3. Did you find or come across code that implements any part of this assignment? If so, include full details.

### Your Answer

1. No 2. No 3. No

## 7 Programming (98 points)

### 7.1 The Task

In the programming section you will implement a named entity recognition system using Hidden Markov Models (HMMs). Named entity recognition (NER) is the task of classifying named entities, typically proper nouns, into pre-defined categories, such as person, location, or organization. Consider the example sequence below, where each word is appended with a tab and then its tag:

"	O
Rhinestone	B-ORG
Cowboy	I-ORG
"	O
(	O
Larry	B-PER
Weiss	I-PER
)	O
-	O
3:15	O

Rhinestone and Cowboy are labeled as an organization (ORG), while Larry and Weiss is labeled as a person (PER). Words that aren't named entities are assigned the O tag. The B- prefix indicates that a word is the beginning of an entity, while the I- prefix indicates that the word is inside the entity.

NER is an incredibly important task for a machine to analyze and interpret a body of natural language text. For example, when designing a system that automatically summarizes news articles, it is important to recognize the key subjects in the articles. Another example is designing a trivia bot. If you can quickly extract the named entities from the trivia question, you may be able to more easily query your knowledge base (e.g. type a query into Google) to request information about the answer to the question.

On a technical level, the main task is to implement an algorithm to learn the HMM parameters given the training data and then implement the forward-backward algorithm to perform a smoothing query which we can then use to predict the hidden tags for a sequence of words.

### 7.2 The Dataset

[WikiANN](#) is a "silver standard" dataset that was generated without human labelling. The English Abstract Meaning Representation (AMR) corpus and DBpedia features were used to train an automatic classifier to label Wikipedia articles. These labels were then propagated throughout other Wikipedia articles using the Wikipedia's cross-language links and redirect links. Afterwards, another tagger that self-trains on the existing tagged entities was used to label all other mentions of the same entities, even those with different morphologies (prefixes and suffixes that modify a word in other languages). Finally, the amassed training examples were filtered by "commonness" and "topical relatedness" to pick more relevant training data.

The WikiANN dataset provides labelled entity data for Wikipedia articles in 282 languages. We will be primarily using the English subset, which contains 14,000 training examples and 3,300 test examples, and the French subset, which contains around 7,500 training examples and 300 test examples.

## 7.3 File Formats

The contents and formatting of each of the files in the handout folder is explained below.

1. **train.txt** This file contains labeled text data that you will use in training your model in the Learning problem (Section 7.4). Specifically, the text contains one word per line that has already been preprocessed, cleaned and tokenized. Every sequence has the following format:

```
<Word0>\t<Tag0>\n<Word1>\t<Tag1>\n ... <WordN>\t<TagN>\n
```

where every `<WordK>\t<TagK>` unit token is separated by a newline. Between each sequence is an empty line. If we have two three-word sequences in our data set, the data will look like so:

```
<Word0>\t<Tag0>\n
<Word1>\t<Tag1>\n
<Word2>\t<Tag2>\n
\n
<Word0>\t<Tag0>\n
<Word1>\t<Tag1>\n
<Word2>\t<Tag2>
```

*Note:* Word 2 of the second sequence does not end with a newline because it is the end of the data set.

2. **validation.txt:** This file contains labeled validation data that you will use to evaluate your model. This file has the same format as **train.txt**.
3. **index\_to\_word.txt, index\_to\_tag.txt:** These files contain a list of all words or tags that appear in the data set. The format is simple:

<b>index_to_word.txt</b>	<b>index_to_tag.txt</b>
<Word0>\n	<Tag0>\n
<Word1>\n	<Tag1>\n
<Word2>\n	<Tag2>\n
⋮	⋮

In your functions, you will convert the string representation of words or tags to indices corresponding to the location of the word or tag in these files. For example, if *Austria* is on line 729 of **index\_to\_word.txt**, then all appearances of *Austria* in the data sets should be converted to the index 729. This index will also correspond to locations in the parameter matrices. For example, the word *Austria* corresponds to the parameters in column 729 of the matrix stored in **hmmemit.txt**. This will be useful for your forward-backward algorithm implementation (see Section 7.5).

4. **predicted.txt:** This file contains labeled data that you will use to debug your implementation. The labels in this file are generated by running a reference implementation using the features from **train.txt**. This file has the same format as **train.txt**.
5. **metrics.txt:** This file contains the metrics you will compute for the validation data. The first line should contain the average log likelihood, and the second line should contain the prediction accuracy. There should be a single space after the colon preceding the metric value; see the reference output file for more detail.
6. **hmmtrans.txt, hmmemit.txt, hmminit.txt:** These files contain pre-trained model parameters of an HMM that you can use to test your implementation of the Learning and Evaluation and Decoding problems (Sections 7.4, 7.5). The formats of the first two files are the same; each line in these files



consists of a conditional probability distribution. In the case of transition probabilities, this distribution corresponds to the probability of transitioning into another state, given a current state. Similarly, in the case of emission probabilities, this distribution corresponds to the probability of emitting a particular symbol, given a current state. Elements in the same row are separated by a space. Each row corresponds to a line of text, using `\n` to create new lines.

**hmmtrans.txt:**

```
<ProbS1S1> <ProbS1S2> ... <ProbS1SN>\n
<ProbS2S1> <ProbS2S2> ... <ProbS2SN>\n...
```

**hmmemit.txt:**

```
<ProbS1Word1> <ProbS1Word2> ... <ProbS1WordN>\n
<ProbS2Word1> <ProbS2Word2> ... <ProbS2WordN>\n...
```

The format of **hmminit.txt** is similarly defined except that it only contains a single probability distribution over starting states. Therefore, each row only has a single element.

**hmminit.txt:**

```
<ProbS1>\n
<ProbS2>\n...
```

## 7.4 Learning

Your first task is to write a program `learnhmm.py` to learn the Hidden Markov Model parameters needed to apply the forward-backward algorithm (See Section 7.5). There are three sets of parameters that you will need to estimate: the initialization probabilities  $\pi$ , the transition probabilities  $\mathbf{B}$ , and the emission probabilities  $\mathbf{A}$ . For this assignment, we model each of these probabilities using a multinomial distribution with parameters  $\pi_j = P(Y_1 = s_j)$ ,  $B_{jk} = P(Y_t = s_k \mid Y_{t-1} = s_j)$ , and  $A_{jk} = P(X_t = k \mid Y_t = s_j)$ . These can be estimated using maximum likelihood, which results in the following parameter estimates:

1.  $P(Y_1 = s_j) = \pi_j = \frac{N_{Y_1=s_j}+1}{\sum_{p=1}^J (N_{Y_1=s_p}+1)}$ , where  $N_{Y_1=s_j}$  equals the number of times state  $s_j$  is associated with the first word of a sentence in the training data set.
2.  $P(Y_t = s_k \mid Y_{t-1} = s_j) = B_{jk} = \frac{N_{Y_t=s_k, Y_{t-1}=s_j}+1}{\sum_{p=1}^J (N_{Y_t=s_p, Y_{t-1}=s_j}+1)}$ , where  $N_{Y_t=s_k, Y_{t-1}=s_j}$  is the number of times state  $s_j$  is followed by state  $s_k$  in the training data set.
3.  $P(X_t = k \mid Y_t = s_j) = A_{jk} = \frac{N_{X_t=k, Y_t=s_j}+1}{\sum_{p=1}^M (N_{X_t=p, Y_t=s_j}+1)}$ , where  $N_{X_t=k, Y_t=s_j}$  is the number of times that the state  $s_j$  is associated with the word  $k$  in the training data set.

Note we add 1 to each count to make a pseudocount. This is slightly different from pure maximum likelihood estimation, but it is useful in improving performance when evaluating unseen cases during evaluation of your validation set.

Your implementation should read in the training data set (**train.txt**), and then estimate  $\pi$ ,  $\mathbf{B}$ , and  $\mathbf{A}$  using the above maximum likelihood solutions with pseudocounts.

Your outputs should be in the same format as **hmminit.txt**, **hmmtrans.txt**, and **hmmemit.txt** (including the same number of decimal places to ensure there are no rounding errors during prediction). The autograder runs and evaluates the output from the files generated, using the following command:

```
$ python3 learnhmm.py [args...]
```

Where `[args...]` is a placeholder for six command-line arguments: `<train_input>` `<index_to_word>` `<index_to_tag>` `<hmminit>` `<hmmemit>` `<hmmtrans>`. These arguments are described below:

1. `<train_input>`: path to the training input `.txt` file (see Section 7.2)
2. `<index_to_word>`: path to the `.txt` file that specifies the dictionary mapping from words to indices. The tags are ordered by index, with the first word having index of 0, the second word having index of 1, etc.
3. `<index_to_tag>`: path to the `.txt` file that specifies the dictionary mapping from tags to indices. The tags are ordered by index, with the first tag having index of 0, the second tag having index of 1, etc.
4. `<hmminit>`: path to output `.txt` file to which the estimated initialization probabilities ( $\pi$ ) will be written. The file output to this path should be in the same format as the handout `hmminit.txt` (see Section 7.2).
5. `<hmmemit>`: path to output `.txt` file to which the emission probabilities (**A**) will be written. The file output to this path should be in the same format as the handout `hmmemit.txt` (see Section 7.2)
6. `<hmmtrans>`: path to output `.txt` file to which the transition probabilities (**B**) will be written. The file output to this path should be in the same format as the handout `hmmtrans.txt` (see Section 7.2).

As an example, the following command would run your program on the toy dataset provided in the handout.

```
$ python3 learnhmm.py toy_data/train.txt toy_data/index_to_word.txt \
toy_data/index_to_tag.txt toy_data/hmminit.txt toy_data/hmmemit.txt \
toy_data/hmmtrans.txt
```

After running the command above, the `<hmminit>`, `<hmmemit>`, and `<hmmtrans>` output files should match the reference files provided in the `toy_output` directory.

## 7.5 Evaluation and Decoding

### 7.5.1 Forward Backward Algorithm and Minimal Bayes Risk Decoding

Your next task is to implement the forward-backward algorithm. Suppose we have a set of sequence consisting of  $T$  words,  $x_1, \dots, x_T$ . Each word is associated with a label  $Y_t \in \{1, \dots, J\}$ . In the forward-backward algorithm we seek to approximate  $P(Y_t | x_{1:T})$  up to a multiplication constant. This is done by first breaking  $P(Y_t | x_{1:T})$  into a “forward” component and a “backward” component as follows:

$$\begin{aligned} P(Y_t = s_j | x_{1:T}) &\propto P(Y_t = s_j, x_{t+1:T} | x_{1:t}) \\ &\propto P(Y_t = s_j | x_{1:t}) P(x_{t+1:T} | Y_t = s_j, x_{1:t}) \\ &\propto P(Y_t = s_j | x_{1:t}) P(x_{t+1:T} | Y_t = s_j) \\ &\propto P(Y_t = s_j, x_{1:t}) P(x_{t+1:T} | Y_t = s_j) \end{aligned}$$

where  $P(Y_t = s_j | x_1, \dots, x_t)$  and  $P(x_{t+1}, \dots, x_T | Y_t = s_j)$  are computed by bottom-up dynamic programming approach.

#### Forward Algorithm

Define  $\alpha_t(s_j) = P(Y_t = s_j, x_{1:t})$ . This can be rearranged into the following expression (the full derivation can be found in the lecture notes):

$$\alpha_t(s_j) = A_{jx_t} \sum_k B_{kj} \alpha_{t-1}(k) \quad (3)$$

Using this definition, the  $\alpha$ 's can be computed using the following dynamic programming procedure:

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

#### Backward Algorithm

Define  $\beta_t(s_j) = P(x_{t+1:T} | Y_t = s_j)$ . This can be rearranged into the following expression:

$$\beta_t(s_j) = \sum_{k=1}^J A_{kx_{t+1}} \beta_{t+1}(s_k) B_{jk} \quad (4)$$

Just like the  $\alpha$ 's, the  $\beta$ 's can also be computed using the following dynamic programming procedure:

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

**Forward-Backward Algorithm** As stated above, the goal of the Forward-Backward algorithm is to compute  $P(Y_t = s_j \mid x_{1:T})$ . This can be done using the following equation:

$$P(Y_t = s_j \mid x_{1:T}) \propto P(Y_t = s_j, x_{1:t})P(x_{t+1:T} \mid Y_t = s_j)$$

After running your forward and backward passes through the sequence, you are now ready to estimate the conditional probabilities as:

$$P(Y_t \mid x_{1:t}) \propto \alpha_t \odot \beta_t$$

where  $\odot$  is the element-wise product.

**Minimum Bayes Risk Prediction** We will assign tags using the minimum Bayes risk predictor, defined for this problem as follows:

$$\hat{Y}_t = \operatorname{argmax}_{j \in \{1, \dots, J\}} P(Y_t = s_j \mid x_{1:T})$$

To resolve ties, select the tag that appears earlier in the `<index_to_tag>` input file.

**Computing the Log Likelihood of a Sequence** When we compute the log likelihood of a sequence, we are interested in the computing the quantity  $\log(P(x_{1:T}))$ . We can rewrite this in terms of values we have already computed in the forward-backward algorithm as follows:

$$\begin{aligned} \log P(x_{1:T}) &= \log \left( \sum_j P(x_{1:T}, Y_t = s_j) \right) \\ &= \log \left( \sum_j \alpha_T(s_j) \right) \end{aligned}$$

### 7.5.2 Implementation Details

You should now write a program `forwardbackward.py` that implements the forward-backward algorithm. The program will read in validation data and the parameter files produced by `learnhmm.py`. The autograder runs and evaluates the output from the files generated, using the following command:

```
$ python3 forwardbackward.py [args...]
```

Where `[args...]` is a placeholder for eight command-line arguments: `<validation_input>` `<index_to_word>` `<index_to_tag>` `<hmm_init>` `<hmm_emit>` `<hmm_trans>` `<predicted_file>` `<metric_file>`.

These arguments are described in detail below:

1. `<validation_input>`: path to the validation input `.txt` file that will be evaluated by your forward backward algorithm (see Section 7.2)
2. `<index_to_word>`: path to the `.txt` file that specifies the dictionary mapping from words to indices. The tags are ordered by index, with the first word having index of 0, the second word having index of 1, etc. This is the same file as was described for `learnhmm.py`.
3. `<index_to_tag>`: path to the `.txt` file that specifies the dictionary mapping from tags to indices. The tags are ordered by index, with the first tag having index of 0, the second tag having index of 1, etc. This is the same file as was described for `learnhmm.py`.
4. `<hmm_init>`: path to input `.txt` file which contains the estimated initialization probabilities ( $\pi$ ).
5. `<hmm_emit>`: path to input `.txt` file which contains the emission probabilities (**A**).
6. `<hmm_trans>`: path to input `.txt` file which contains transition probabilities (**B**).
7. `<predicted_file>`: path to the output `.txt` file to which the predicted tags will be written. The file should be in the same format as the `<validation_input>` file.
8. `<metric_file>`: path to the output `.txt` file to which the metrics will be written.

As an example, the following command would run your program on the toy dataset provided in the handout.

```
$ python3 forwardbackward.py toy_data/validation.txt \
toy_data/index_to_word.txt toy_data/index_to_tag.txt \
toy_data/hmm_init.txt toy_data/hmm_emit.txt \
toy_data/hmm_trans.txt toy_data/predicted.txt \
toy_data/metrics.txt
```

After running the command above, the `<predicted_file>` output should be:

```
fish    D
eat     C
you     D
```

And the `<metric_file>` output should be:

```
Average Log-Likelihood: -3.0438629330222424
Accuracy: 0.3333333333333333
```

where average log-likelihood and accuracy are evaluated over the validation set.

Take care that your output has the exact same format as shown above. There should be a single space after the colon preceding the metric value (e.g. a space after `Average Log-Likelihood:`). Each line should be terminated by a Unix line ending `\n`.

### 7.5.3 Log-Space Arithmetic for Avoiding Underflow

Handling underflow properly is a critical step in implementing an HMM. The most generalized way of handling numerical underflow due to products of small positive numbers (like probabilities) is to calculate everything in log-space, i.e., represent every quantity by their logarithm.

For this homework, using log-space starts with transforming Eq.(3) and Eq.(4) into logarithmic form - you may find the recitation handout helpful. Please use base  $e$  (natural log) for logarithm calculation.

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](#). We provide the pseudocode for this trick in Algorithm 1:

---

**Algorithm 1** Log-Sum-Exp Trick

---

```
1: procedure LOGSUMEXPTRICK( $(v_1, v_2, \dots, v_n)$ )
2:    $m = \max(v_i)$  for  $i = \{1, 2, \dots, n\}$ 
3:   return  $m + \log(\sum_i \exp(v_i - m))$ 
```

---

**Note:** The autograder test cases account for numerical underflow using the Log-Sum-Exp Trick. If you do not implement `forwardbackward.py` with the trick, you might only receive partial credit.

## 7.6 Gradescope Submission

You should submit your `learnhmm.py` and `forwardbackward.py` to Gradescope. **Any other files will be deleted.** Please do not use other file names. This will cause problems for the autograder to correctly detect and run your code. Please go through the appendix at the end for information on starter-code.

Some additional tips: Make sure to read the autograder output carefully. The autograder for Gradescope prints out some additional information about the tests that it ran. For this programming assignment we've specially designed some buggy implementations that you might implement and will try our best to detect those and give you some more useful feedback in Gradescope's autograder. Make wise use of autograder's output for debugging your code.

*Note:* For this assignment, you have 10 submissions to Gradescope before the deadline, but only your last submission will be graded.