

# HOMWORK 7: HIDDEN MARKOV MODELS

10-301/10-601 Introduction to Machine Learning (Fall 2020)

<https://www.cs.cmu.edu/~10601/>

DUE: Thursday, November 19, 2020 11:59 PM

**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 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: <https://www.cs.cmu.edu/~10601/>
- **Late Submission Policy:** See the late submission policy here: <https://www.cs.cmu.edu/~10601/>
- **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 must be written in LaTeX. Each derivation/proof should be completed in the boxes provided. If you do not follow the template, your assignment may not be graded correctly by our AI assisted grader.
  - **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 3.6.9, OpenJDK 11.0.5, g++ 7.4.0) and versions of permitted libraries (e.g. `numpy` 1.17.0 and `scipy` 1.4.1) match those used on Gradescope. You have unlimited Gradescope programming submissions. However, we recommend debugging your implementation on your local machine (or the Linux servers) and making sure your code is running correctly first before submitting you code to Gradescope.
- **Materials:** The data that you will need in order to complete this assignment is posted along with the writeup and template on Piazza.

**Linear Algebra Libraries** When implementing machine learning algorithms, it is often convenient to have a linear algebra library at your disposal. In this assignment, Java users may use EJML<sup>a</sup> or ND4J<sup>b</sup> and C++ users Eigen<sup>c</sup>. Details below. (As usual, Python users have NumPy.)

**EJML for Java** EJML is a pure Java linear algebra package with three interfaces. We strongly recommend using the SimpleMatrix interface. The autograder will use EJML version 0.38. When compiling and running your code, we will add the additional command line argument `-cp "linalg_lib/ejml-v0.38-libs/*:linalg_lib/nd4j-v1.0.0-beta7-libs/*:./"` to ensure that all the EJML jars are on the classpath as well as your code.

**ND4J for Java** ND4J is a library for multidimensional tensors with an interface akin to Python's NumPy. The autograder will use ND4J version 1.0.0-beta7. When compiling and running your code, we will add the additional command line argument `-cp "linalg_lib/ejml-v0.38-libs/*:linalg_lib/nd4j-v1.0.0-beta7-libs/*:./"` to ensure that all the ND4J jars are on the classpath as well as your code.

**Eigen for C++** Eigen is a header-only library, so there is no linking to worry about—just `#include` whatever components you need. The autograder will use Eigen version 3.3.7. The command line arguments above demonstrate how we will call your code. When compiling your code we will include, the argument `-I./linalg_lib` in order to include the `linalg_lib/Eigen` subdirectory, which contains all the headers.

We have included the correct versions of EJML/ND4J/Eigen in the `linalg_lib.zip` posted on the Piazza Resources page for your convenience. It contains the same `linalg_lib/` directory that we will include in the current working directory when running your tests. Do **not** include EJML, ND4J, or Eigen in your homework submission; the autograder will ensure that they are in place.

---

<sup>a</sup><https://ejml.org>

<sup>b</sup><https://deeplearning4j.org/docs/latest/nd4j-overview>

<sup>c</sup><http://eigen.tuxfamily.org/>

## Written Questions (20 points)

### 1. Multiple Choice

(a) In this section, we will test your understanding of several aspects of HMMs. Shade in the box or circle in the template document corresponding to the correct answer(s) for each of the subparts below.

i. (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

ii. (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  $Y_{t-1}$
- ☐ The current observation  $X_t$  is conditionally independent of  $Y_{t-2}$  given the previous observation  $X_{t-1}$
- ☐ None of the above

- (b) In the remaining subparts, you will always see two quantities and decide what is the strongest relation between them. There is **only one correct answer**. Use the following definitions:  $\alpha_t(s_j) = P(Y_t = s_j, x_{1:t})$  and  $\beta_t(s_j) = P(x_{t+1:T} | Y_t = s_j)$ .

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

- i. (1 point) (**Select one**) Let  $N$  denote the number of possible hidden states. In other words, each  $Y_t \in \{s_i\}_{i=1}^N$ . What is the relation between  $\sum_{i=1}^N (\alpha_5(s_i)\beta_5(s_i))$  and  $P(x_{1:T})$ ? Select only the **strongest** relation that necessarily holds.

☐ =

☐ >

☐ <

☒ ≤

☐ ≥

☐ ?

- ii. (1 point) (**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.

☒ =

☐ >

☐ <

☐ ≤

☐ ≥

☐ ?

- iii. (1 point) (**Select one**) What is the relation between  $\alpha_5(s_i)$  and  $\beta_5(s_i)$ ? Select only the **strongest** relation that necessarily holds.

☐ =

☐ >

☐ <

☐ ≤

☐ ≥

☒ ?

## 2. Viterbi Decoding

- (a) Suppose we have a set of sequence 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:

- $\mathbf{A}$  is the transition matrix:  $A_{jk} = P(Y_t = s_k \mid Y_{t-1} = s_j)$
- $\mathbf{B}$  is the emission matrix:  $B_{jk} = P(X_t = k \mid Y_t = s_j)$
- $\pi$  describes  $Y_1$ 's initialization probabilities:  $\pi_j = P(Y_1 = s_j)$
- $\omega_t(s_k)$  is the maximum product of all the probabilities taken through path  $Y_1, \dots, Y_{t-1}$  that ends with  $Y_t$  at state  $s_k$ .

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

- $b_t(s_k)$  are the backpointers that store the path through hidden states that give 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:

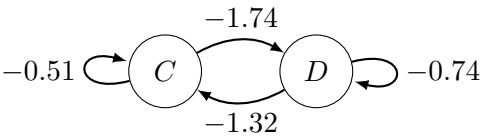
1. Initialize  $\omega_1(s_j) = \pi_j B_{jx_1}$  and  $b_1(j) = j$
2. For  $t > 1$ , we have

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

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 subpart, consider the Hidden Markov Model specified below. Working with probabilities in the logarithm scale, we have that  $\ln P(Y_1 = C) = \ln P(Y_1 = D) = -0.69$ , and the state transition model and emission probability tables are given as follows.



$k$	$\ln P(X_t = k \mid Y_t = C)$	$\ln P(X_t = k \mid Y_t = D)$
1	-0.69	-1.21
2	-0.91	-0.69
3	-2.30	-1.61

We observed  $X_1 = 1$  and  $X_2 = 2$ . (Note that taking the maximum of log probabilities will give you the same result as taking the maximum of probabilities as log is a monotonically increasing function.)

- i. (1 point) Compute  $\ln \omega_1(C)$  and  $\ln \omega_1(D)$ . If your answers involves 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  $\ln \omega_1(C)$ ?

$\ln \omega_1(C)$	Work

What is  $\ln \omega_1(D)$ ?

$\ln \omega_1(D)$	Work

ii. (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.

### 3. Warm-up Exercise: Forward-Backward Algorithm

- (a) (4 points) To help you prepare to implement the HMM forward-backward algorithm (see Section 3 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 unique words used in the training set.

**Training set:**

```
you_D eat_C fish_D
you_D fish_D eat_C
eat_C fish_D
```

Where the training word sequences are:

$$\begin{aligned}\mathbf{x}^{(1)} &= [\text{you} \text{ eat} \text{ fish}]^T \\ \mathbf{x}^{(2)} &= [\text{you} \text{ fish} \text{ eat}]^T \\ \mathbf{x}^{(3)} &= [\text{eat} \text{ fish}]^T\end{aligned}$$

And the corresponding tags are:

$$\begin{aligned}\mathbf{y}^{(1)} &= [D \ C \ D]^T \\ \mathbf{y}^{(2)} &= [D \ D \ C]^T \\ \mathbf{y}^{(3)} &= [C \ D]^T\end{aligned}$$

**Validation set:**

```
fish eat you
```

Where the validation word sequences are:

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

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



- $\alpha_t(j) = P(Y_t = s_j, x_{1:t})$ , which can be computed recursively:
  1.  $\alpha_1(j) = \pi_j B_{jx_1}$ .
  2. For  $t > 1$ ,  $\alpha_t(j) = B_{jx_t} \sum_{k=1}^J \alpha_{t-1}(k) A_{kj}$
- $\beta_t(j) = P(x_{t+1:T} | Y_t = s_j)$ , which can be computed recursively:
  1.  $\beta_T(j) = 1$
  2. For  $1 \leq t \leq T - 1$ ,  $\beta_t(j) = \sum_{k=1}^J B_{kx_{t+1}} \beta_{t+1}(k) A_{jk}$

The following four subparts are meant to encourage you to work through the forward backward algorithm by hand using this validation example. Feel free to use a calculator, being careful to carry enough significant figures through your computations to avoid rounding errors. For each subpart below, please report the requested value in the text box under the subpart (these boxes are only visible in the template document).

**Note:** pseudo count used in section 2 should also used here.

*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.

- (1 point) Compute  $\alpha_2(C)$ , the  $\alpha$  value associated with the tag “C” for the second word in the validation sequence. Please round your answer to **THREE** decimal places.

$\alpha_2(C)$	Work
0.131	

- ii. (1 point) Compute  $\beta_2(D)$ , the  $\beta$  value associated with the tag “D” for the second word in the validation sequence. Please round your answer to **THREE** decimal places.

$\beta_2(D)$	Work
0.25	

- iii. (1 point) Predict the tag for the third word in the validation sequence.

Tag	Work
D	

- iv. (1 point) Compute the log-likelihood for the entire validation sequence, “fish eat you”. Please round your answer to **THREE** decimal places.

Log-Likelihood	Work
-3.044	

- (b) (6 points) Return to these subparts after implementing your `learnhmm.{py|java|cpp}` and `forwardbackward.{py|java|cpp}` functions. Please ensure that you have used the log-sum-exp trick in your programming as described in section 3.3 before answering these empirical questions.

Using the fulldata set **trainwords.txt** in the handout using your implementation of `learnhmm.{py|java|cpp}` to learn parameters for an hmm model using the first 10, 100, 1000, and 10000 sequences in the file. Use these learned parameters to perform prediction on the **trainwords.txt** and the **validationwords.txt** files using your `forwardbackward.{py|java|cpp}`. Construct a plot with number of sequences used for training on the x-axis (log-scale) and average log likelihood across all sequences from the **trainwords.txt** and the **validationwords.txt** on the y-axis (see Section 3 for details on computing the log data likelihood for a sequence). Each table entry is worth 0.5 points. Write the resulting log likelihood values in the table in the template. Include your plot in the large box in the template (2 points). To receive credit for your plot, you must submit a computer generated plot. **DO NOT** hand draw your plot.

- i. (4 points) Fill in this table.

# Sequences	Train Average Log-Likelihood	Validation Average Log-Likelihood
10	??	??
100	??	??
1000	??	??
10000	??	??

- ii. (2 points) Put your plot below:

*Plot*

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

### Answer

I collaborated with Sinduja Sriskanda on high level topics, going through the recitation 8 together until we understood it.

# Programming (80 points)

## 1 The Tasks and Data Sets

In the programming section you will implement a named entity recognition system using HMMs. Named entity recognition is the task of classifying named entities (typically proper nouns) into pre-defined categories (e.g. person, location, organization, etc). Consider the example below, where each word is appended with an underscore and then its tag:

```
All-rounder_O Phill_PER Simmons_PER lead_O Somerset_ORG in_O win_O
```

The words `Phill` and `Simmons` are labeled as a person (`PER`), while the word `Somerset` is labeled as an organization (`ORG`). Words that aren't named entities are assigned the `O` tag.

Named entity is an incredibly important task for a machine to begin to analyze and interpret a body of natural language text. For example, when designing a system that automatically summarizes all of the day's 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 hidden Markov model parameters given the training data and then implement a the forward backward algorithm to perform a smoothing query which we can then use to predict the hidden tags for a sequence of words.

The handout files for this assignments contains several files that you will use in the homework. The contents and formatting of each of these files is explained below.

1. **trainwords.txt** This file contains labeled text data that you will use in training your model in the **Learning** problem. Specifically the text contains one sentence per line that has already been preprocessed, cleaned and tokenized. You should treat every line as a separate sequence and assume that it has the following format:

```
<Word0>_<Tag0> <Word1>_<Tag1> ... <WordN>_<TagN>
```

where every `<WordK>_<TagK>` unit token is separated by white space.

2. **validationwords.txt**: This file contains labeled data that you will use to evaluate your model in the **Experiments** section. This file contains the gold standard labels. This file has the same format as **trainwords.txt**.
3. **index\_to\_word.txt** and **index\_to\_tag.txt**: These files contain a list of all words or tags that appear in the data set. 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 3).
4. **toytrain.txt**, **toyvalidation.txt**, **toy\_index\_to\_word.txt**, and **toy\_index\_to\_tag.txt**: These files are analogous to **trainwords.txt**, **validationwords.txt**, **index\_to\_word.txt**, and **index\_to\_tag.txt** and are used to compare your implementation to your hand calculations in the Written component.
5. **predictvalidation.txt** The **predictvalidation.txt** file contains labeled data that you will use to debug

your implementation. The labels in this file are not gold standard but are generated by running our decoder using the features from **trainwords.txt**. This file has the same format as **trainwords.txt**.

6. **metrics.txt** The **metrics.txt** file contains the metrics you will compute for the dataset. For this assignment, you need to compute the average log likelihood and your prediction accuracy on the validation data. Note that in Named Entity Recognition, F-1 score is a more common metric to evaluate the model performance, here you only need to report your accuracy for tag prediction of each word.
7. **hmmtrans.txt, hmmemit.txt and hmmprior.txt**: These files contain pre-trained model parameters of an HMM that you will use in testing your implementation of the **Learning** and **Evaluation and Decoding** problems. The format of the first two files are analogous and is as follows. Every 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. For example, every line in **hmmtrans.txt** has the following format:

**hmmtrans.txt:**

```
<ProbS1S1> ... <ProbS1SN>
<ProbS2S1> ... <ProbS2SN>...
```

and every line in **hmmemit.txt** has the following format:

**hmmemit.txt:**

```
<ProbS1Word1> ... <ProbS1WordN>
<ProbS2Word1> ... <ProbS2WordN>...
```

In both cases, elements in the same row are separated by white space. Each row corresponds to a line of text (using `\n` to create new lines).

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

**hmmprior.txt:**

```
<ProbS1>
<ProbS2>...
```

Note that the data provided to you is to help in developing your implementation of the HMM algorithms. Your code will be tested on Gradescope using different data with different HMM parameters, likely coming from a different domain although the format will be identical.

## 2 Learning

Your first task is to implement an algorithm to learn the hidden Markov model parameters needed to apply the forward backward algorithm (See Section 3). There are three sets of parameters that you will need to estimate: the initialization probabilities  $\pi$ , the transition probabilities  $\mathbf{A}$ , and the emission probabilities  $\mathbf{B}$ . For this assignment, we model each of these probabilities using a multinomial distribution with parameters  $\pi_j = P(Y_1 = s_j)$ ,  $A_{jk} = P(Y_t = s_k \mid Y_{t-1} = s_j)$ , and  $B_{jk} = P(X_t = k \mid Y_t = s_j)$ . These can be estimated using maximum likelihood, which results in the following parameter estimators:

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.  $P(Y_t = s_k \mid Y_{t-1} = s_j) = A_{jk} = \frac{N_{Y_t=s_k, Y_{t-1}=s_j}+1}{\sum_{p=1}^J(N_{Y_t=s_p, Y_{t-1}=s_j}+1)}$ , 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.  $P(X_t = k \mid Y_t = s_j) = B_{jk} = \frac{N_{X_t=k, Y_t=s_j}+1}{\sum_{p=1}^M(N_{X_t=p, Y_t=s_j}+1)}$ , 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 that for each count, a “+1” is added 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.

You should implement a function that reads in the training data set (**trainwords.txt**), and then estimates  $\pi$ ,  $A$ , and  $B$  using the above maximum likelihood solutions.

Your outputs should be in the same format as **hmmprior.txt**, **hmmtrans.txt**, and **hmmemit.txt** (including the same number of decimal places to ensure there are no rounding errors during prediction). The autograder will use the following commands to call your function:

For Python: `$ python learnhmm.py [args...]`  
 For Java: `$ javac -cp "./lib/ejml-v0.33-libs/*:./" learnhmm.java;`  
`java -cp "./lib/ejml-v0.33-libs/*:./" learnhmm [args...]`  
 For C++: `$ g++ -g -std=c++11 -I./lib learnhmm.cpp; ./a.out [args...]`

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

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

### 3 Evaluation and Decoding

#### 3.1 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  $P(Y_t = s_j | x_1, \dots, x_t)$  is computed by passing forward recursively through the model and  $P(x_{t+1}, \dots, x_T | Y_t = s_j)$  is computed by passing recursively backwards through the model.

#### Forward Algorithm

Define  $\alpha_t(j) = P(Y_t = s_j, x_{1:t})$ . We can rearrange our definition of  $\alpha_t(j)$  as follows:

$$\begin{aligned} \alpha_t(j) &= P(Y_t = s_j, x_{1:t}) \\ &= \sum_k P(Y_t = s_j, Y_{t-1} = s_k, x_{1:t}) \\ &= \sum_k P(Y_{t-1} = s_k, x_{1:t} | Y_t = s_j) P(Y_t = s_j) \\ &= \sum_k P(x_t | Y_t = s_j) P(Y_{t-1} = s_k, x_{1:t-1} | Y_t = s_j) P(Y_t = s_j) \\ &= P(x_t | Y_t = s_j) \sum_k P(Y_t = s_j, Y_{t-1} = s_k, x_{1:t-1}) \\ &= P(x_t | Y_t = s_j) \sum_k P(Y_t = s_j, x_{1:t-1} | Y_{t-1} = s_k) P(Y_{t-1} = s_k) \\ &= P(x_t | Y_t = s_j) \sum_k P(Y_t = s_j | Y_{t-1} = s_k) P(x_{1:t-1} | Y_{t-1} = s_k) P(Y_{t-1} = s_k) \\ &= P(x_t | Y_t = s_j) \sum_k P(Y_t = s_j | Y_{t-1} = s_k) P(Y_{t-1} = s_k, x_{1:t-1}) \\ &= B_{jx_t} \sum_k A_{kj} \alpha_{t-1}(k) \end{aligned} \tag{3}$$

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

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



**Backward Algorithm** Define  $\beta_t(j) = P(x_{t+1:T} \mid Y_t = s_j)$ . We can rearrange our definition of  $\beta_t(j)$  as follows:

$$\begin{aligned}
\beta_t(j) &= P(x_{t+1:T} \mid Y_t = s_j) \\
&= \sum_{k=1}^J P(Y_{t+1} = s_k, x_{t+1:T} \mid Y_t = s_j) \\
&= \sum_{k=1}^J P(x_{t+1:T} \mid Y_t = s_j, Y_{t+1} = s_k) P(Y_{t+1} = s_k \mid Y_t = s_j) \\
&= \sum_{k=1}^J P(x_{t+1} \mid Y_{t+1} = s_k) P(x_{t+2:T} \mid Y_{t+1} = s_k) P(Y_{t+1} = s_k \mid Y_t = s_j) \\
&= \sum_{k=1}^J B_{kx_{t+1}} \beta_{t+1}(k) A_{jk}
\end{aligned} \tag{4}$$

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

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

**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:

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

where  $\circ$  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  $P(x_{1:T})$ . We can rewrite this in terms of values we have already computed in the forward-backward algorithm as follows:

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

### 3.2 Implementation Details

You should now implement your forward-backward algorithm as a program, `forwardbackward.{py|java|cpp}`. The program will read in validation data and the parameter files produced by `learnhmm.{py|java|cpp}`. The autograder will use the following commands to call your function:

```
For Python: $ python forwardbackward.py [args...]
For Java:   $ javac -cp "./lib/ejml-v0.33-libs/*:./" forwardbackward.java;
             $ java -cp "./lib/ejml-v0.33-libs/*:./" forwardbackward [args...]
For C++:    $ g++ -g -std=c++11 -I./lib forwardbackward.cpp; ./a.out [args...]
```

Where above `[args...]` is a placeholder for seven command-line arguments: `<validation_input>` `<index_to_word>` `<index_to_tag>` `<hmmprior>` `<hmmemit>` `<hmmtrans>` `<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 1)
2. `<index_to_word>`: path to the `.txt` that specifies the dictionary mapping from words to indices. The tags are ordered by index, with the first word having index of 1, the second word having index of 2, etc. This is the same file as was described for `learnhmm.{py|java|cpp|m}`.
3. `<index_to_tag>`: path to the `.txt` that specifies the dictionary mapping from tags to indices. The tags are ordered by index, with the first tag having index of 1, the second tag having index of 2, etc. This is the same file as was described for `learnhmm.{py|java|cpp|m}`.
4. `<hmmprior>`: path to input `.txt` file which contains the estimated prior ( $\pi$ ).
5. `<hmmemit>`: path to input `.txt` file which contains the emission probabilities (**B**).
6. `<hmmtrans>`: path to input `.txt` file which contains transition probabilities (**A**).
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.

Example command for python users:

```
$ python forwardbackward.py toy_data/toy_train.txt \
toy_data/toy_index_to_word.txt toy_data/toy_index_to_tag.txt \
toy_data/toy_hmmprior.txt toy_data/toy_hmmemit.txt \
toy_data/toy_hmmtrans.txt toy_data/toy_predicted.txt \
toy_data/toy_metrics.txt
```

After running the command above, the <predicted\_file> output should be:

```
you_B eat_A fish_B
you_B fish_B eat_A
eat_A fish_B
```

And the <metric\_file> output should be:

```
Average Log-Likelihood: -2.776793
Accuracy: 1.0
```

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`.

### 3.3 Log-Space Arithmetic for Avoiding Underflow

Handling underflow properly is a critical step in implementing an HMM. And 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 trick starts with transforming Eq.(3) and Eq.(4) into logarithmic form (how to do that is straightforward and left as an exercise). Please use  $e$  as the base for logarithm calculation (natural log).

After transforming the equations into log form, you may discover calculation 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 pseudo code 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**  $m + \log(\sum_i \exp(v_i - m))$
-

## 4 Gradescope Submission

You should submit your `learnhmm.{py|java|cpp}` and `forwardbackward.{py|java|cpp}` to Gradescope. Note: please do not use other file names. This will cause problems for the autograder to correctly detect and run your 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 do, and 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 unlimited submissions to Gradescope before the deadline, but only your last submission will be graded.