

BDP: Homework 2: Executing Baum Welch Program

Due on Friday, December 22nd, 2017

Joonho Kwon

Hani Ramadhan - 201793254

Contents

1	Preliminary Informations	3
2	Preparation	3
2.1	Source Code Modification	3
2.2	BaumWelch Algorithm Input File	4
3	Class Diagram	4
4	Program Execution	6

1 Preliminary Informations

The purpose of this homework report is to show the success of executing BaumWelch algorithm [1] in hadoop using the source code located in <https://github.com/DhruvKumar/Baum-Welch>. This report is consisted of preparation, class diagram of the source code, and the execution result.

2 Preparation

This section will discuss about the source code provided for this project and setup for program's input file. The list of source code and its treatment is provided in Table 1 and Figure 1.

2.1 Source Code Modification

Source code modifications and reasons why some files are not used are briefly described below:

- **BaumWelchDriver.java** modification: Several options were not implemented in DefaultOptionCreator class. Thus, the options should be customized individually. This leads to make the model always built randomly at the beginning.
 - Removed all addOption(...) function calls that were not contained in DefaultOptionCreator:
 - * addOption(DefaultOptionCreator.modelInOption()...
 - * addOption(DefaultOptionCreator.numHiddenStates...
 - * addOption(DefaultOptionCreator.numObservedStat...
 - Recreate removed all addOption() function call by adding own created DefaultOptionBuilder()...
 - Reassigned the appropriate getOption(...) commands as reflected in own created addOption() function calls.
 - Set buildRandom boolean value as true. Thus, the model always will be built randomly at initialization.
- **BaumWelchMapper.java** modification: There are some elements in BaumWelchConfigKeys.java which were not configured properly. They are *HIDDEN_STATES_MAP_PATH* and *EMITTED_STATES_MAP_PATH*. Hence, some modifications were applied.

Table 1: List of the BaumWelch source code files and their treatment

File Name	Treatment
BaumWelchCombiner.java	not used
BaumWelchConfigKeys.java	used, no modification
BaumWelchDriver.java	used, modified
BaumWelchDriver.java	not used
BaumWelchMapper.java	used, modified
BaumWelchModel.java	not used
BaumWelchReducer.java	used, no modification
BaumWelchUtils.java	used, no modification
IntArrayWritable.java	used, no modification
MapWritableCache.java	not used after modification of BaumWelchMapper.java

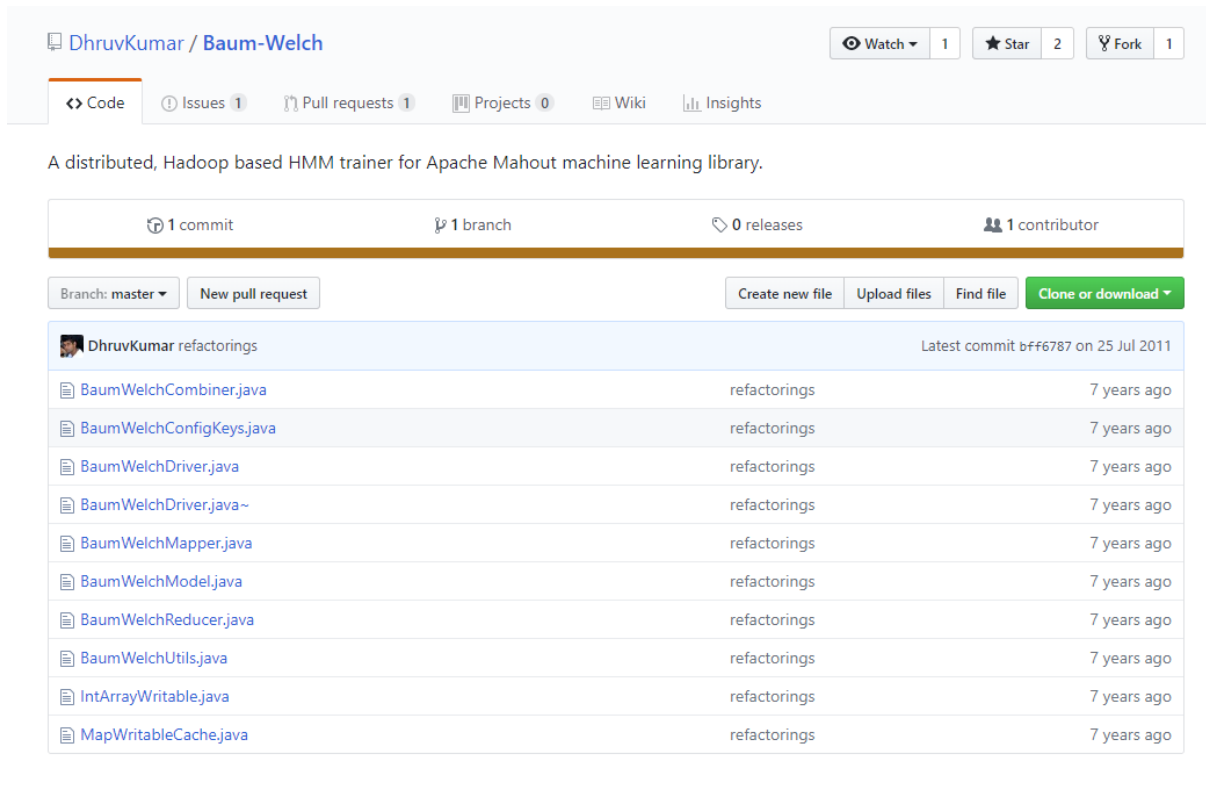


Figure 1: Source Code of Baum Welch Algorithm

- Removal of *hiddenStatesWritableMap* and *emittedStatesWritableMap* variables. Which followed by removal to all of the related commands and variables. This also annuls the full usage/role of *MapWritableCache* in the project.
- BaumWelchCombiner.java unused: never called in the project
- BaumWelchModel.java unused: never called in the project
- MapWritableCache.java unused: already explained

2.2 BaumWelch Algorithm Input File

BaumWelch's input file requires sequence file(s) that should be put in input folder. Sequence file is a flat file which contains serialized key-value pair [2]. This flat file is commonly used in MapReduce task. BaumWelch program, as defined in BaumWelchMapper.java, needs LongWritable type as key and IntArrayWritable as value. Thus, the source code provided in <https://examples.javacodegeeks.com/enterprise-java/apache-hadoop/hadoop-sequence-file-example/> helped to initially create sequence file(s) as much as needed in BaumWelch input folder. The interaction of both programs are depicted in Figure 2.

3 Class Diagram

The class diagram of the BaumWelch program is presented in Figure 3, minus the BaumWelchModel class.

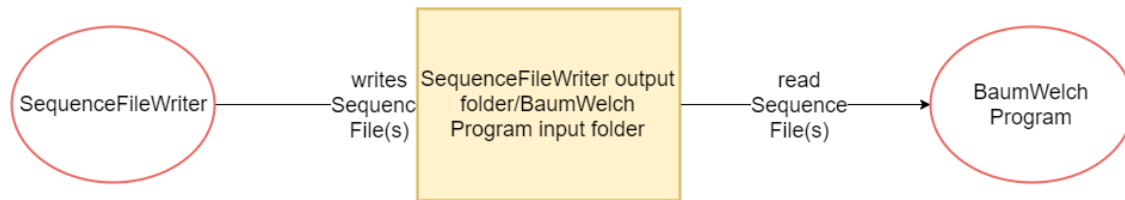


Figure 2: The interaction between Sequence File Writer and Baum Welch Program

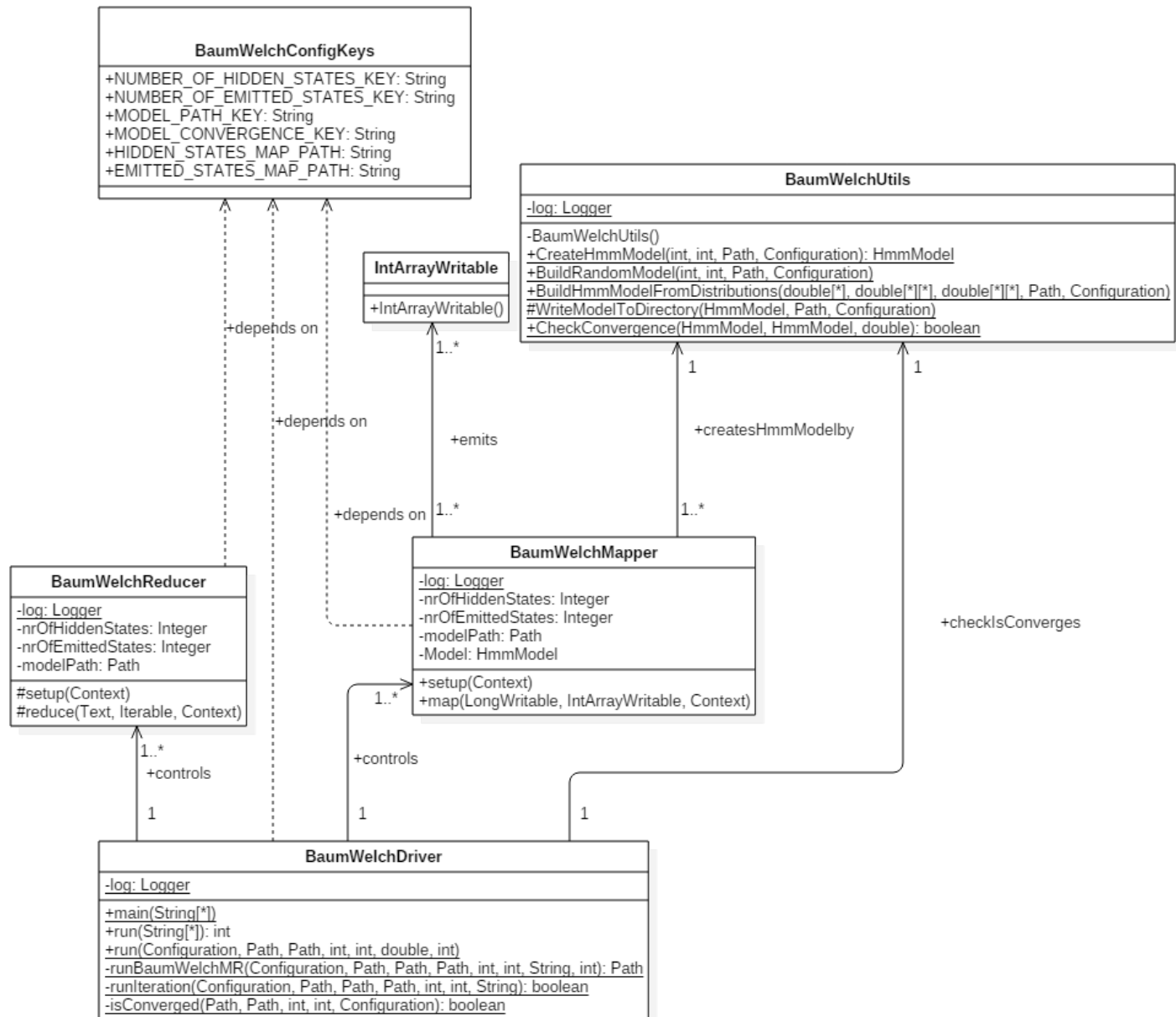


Figure 3: Class Diagram of Baum Welch program

```
hadoop jar BaumWelch.jar --input hani/input --output hani/output
--maxIter 10 --hmmModelInputPath hani/model --numHiddenStates 3
--numObservedStates 4
```

Figure 4: Example of BaumWelch program Execution, red colored means arguments

```
hduser@pc11: ~/hani
hduser@pc11:~/hani$ hadoop jar BaumWelch.jar --input hani/input --output hani/output --maxIter 10 --hmmModelInputPath hani/model --numHiddenStates 3 --numObservedStates 4
17/12/23 02:55:41 INFO common.AbstractJob: Command line arguments: [--convergenceDelta=[0.5], --endPhase=[2147483647], --hmmModelInputPath=[hani/model], --input=[hani/input], --maxIter=[10], --numHiddenStates=[3], --numObservedStates=[4], --output=[hani/output], --startPhase=[0], --tempDir=[temp]]
Num Hidden: 3
17/12/23 02:55:41 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
17/12/23 02:55:41 INFO common.HadoopUtil: Deleting hani/model
17/12/23 02:55:41 INFO compress.CodecPool: Got brand-new compressor [.deflate]
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Initial Distribution Map: State 0 = 0.5582831611058185
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Transition Distribution Map Inner: (0, 0) = (0, 0.4694982088853295)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Transition Distribution Map Inner: (0, 1) = (1, 0.03013048223999798)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Transition Distribution Map Inner: (0, 2) = (2, 0.5003713089906707)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Emission Distribution Map Inner: (0, 0) = (0, 0.34128091274851036)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Emission Distribution Map Inner: (0, 1) = (1, 0.08244654722481974)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Emission Distribution Map Inner: (0, 2) = (2, 0.22627027948530268)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Emission Distribution Map Inner: (0, 3) = (3, 0.35000226054136724)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Initial Distribution Map: State 1 = 0.25019435683106134
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Transition Distribution Map Inner: (1, 0) = (0, 0.13262099831883284)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Transition Distribution Map Inner: (1, 1) = (1, 0.45560472772553)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Transition Distribution Map Inner: (1, 2) = (2, 0.41177427395563715)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Emission Distribution Map Inner: (1, 0) = (0, 0.23635836943951688)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Emission Distribution Map Inner: (1, 1) = (1, 0.004051745156192072)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Emission Distribution Map Inner: (1, 2) = (2, 0.29355768254899056)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Emission Distribution Map Inner: (1, 3) = (3, 0.4660322028553005)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Initial Distribution Map: State 2 = 0.1915224820631201
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Transition Distribution Map Inner: (2, 0) = (0, 0.37496026199984367)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Transition Distribution Map Inner: (2, 1) = (1, 0.49459751398175045)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Transition Distribution Map Inner: (2, 2) = (2, 0.1304422240184059)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Emission Distribution Map Inner: (2, 0) = (0, 0.39118518150589005)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Emission Distribution Map Inner: (2, 1) = (1, 0.42203433751088065)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Emission Distribution Map Inner: (2, 2) = (2, 0.10465796874596425)
17/12/23 02:55:41 INFO BaumWelchUtils: BuildRandomModel Emission Distribution Map Inner: (2, 3) = (3, 0.08212251223726504)
17/12/23 02:55:41 INFO BaumWelchUtils: Wrote random Initial Distribution Map to hani/model/part-randomSeed
17/12/23 02:55:41 INFO BaumWelchUtils: Writing Transition Distribution Map Key, Value = (TRANSIT_2, org.apache.hadoop.io.MapWritable@6578b4cc3)
17/12/23 02:55:41 INFO BaumWelchUtils: Writing Transition Distribution Map Key, Value = (TRANSIT_1, org.apache.hadoop.io.MapWritable@686a2b45d)
17/12/23 02:55:41 INFO BaumWelchUtils: Writing Transition Distribution Map Key, Value = (TRANSIT_0, org.apache.hadoop.io.MapWritable@6572cd92f)
17/12/23 02:55:41 INFO BaumWelchUtils: Wrote random Transition Distribution Map to hani/model/part-randomSeed
17/12/23 02:55:41 INFO BaumWelchUtils: Writing Emission Distribution Map Key, Value = (EMIT_0, org.apache.hadoop.io.MapWritable@fc9c4ee6)
```

Figure 5: Start of Baum Welch program execution

4 Program Execution

The example of the Baum Welch algorithm program execution is described in Figure 4. Options *-input*, *-output*, and *-maxIter* are built in from the program and cannot be neglected to run the program. The customized options are *-hmmModelInputPath*, *-numHiddenStates*, and *numObservedStates*. The options' information are detailed in Table 2. The screenshot of the start and the ending are presented in Figure 5 and Figure 6.

References

- [1] D. Kumar, "Baum-Welch: A distributed, Hadoop based HMM trainer for Apache Mahout machine learning library,," Dec. 2015. original-date: 2011-07-24T16:35:10Z.
- [2] R. Jhajj, "Hadoop Sequence File Example."

```

17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel stateID = key.charAt(8) = 2
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Transition Matrix (2, 0) = 6.872426218891877E-8
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Transition Matrix (2, 1) = 0.5076844193274871
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Transition Matrix (2, 2) = 0.4923155119482507
17/12/23 02:57:06 INFO BaumWelchUtils: Entering Create Hmm Model. Model Path = hani/output/model-6
17/12/23 02:57:06 INFO BaumWelchUtils: Create Hmm Model. ModelFiles Path = hani/output/model-6/*
17/12/23 02:57:06 INFO BaumWelchUtils: Create Hmm Model. File System = DFS[DFSClient[clientName=DFSClient_NONMAPREDUCE_-66445
3659_1, ugi=hduser (auth:SIMPLE)]]
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Adding File Match hdfs://164.125.37.221:54310/user/hduser/hani/output/
model-6/part-r-00000
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Matching Seq File Key = EMIT_0
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Emission Matrix (0, 0) = 3.141467416436764E-11
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Emission Matrix (0, 1) = 1.724951094275974E-15
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Emission Matrix (0, 2) = 2.664678137618084E-13
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Emission Matrix (0, 3) = 0.99999999999683172
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Matching Seq File Key = EMIT_1
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Emission Matrix (1, 0) = 1.122296186082004E-13
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Emission Matrix (1, 1) = 0.03241785723468804
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Emission Matrix (1, 2) = 0.48617101511616134
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Emission Matrix (1, 3) = 0.48141112764903843
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Matching Seq File Key = EMIT_2
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Emission Matrix (2, 0) = 0.5166010076327254
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Emission Matrix (2, 1) = 0.4822319007766118
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Emission Matrix (2, 2) = 0.0011670772514247706
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Emission Matrix (2, 3) = 1.4338237927077507E-8
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Matching Seq File Key = INITIAL
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Initial Prob Adding Key, Value = (0 2.0746956624734732E-17)
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Initial Prob Adding Key, Value = (1 0.06652930665756908)
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Initial Prob Adding Key, Value = (2 0.9334706933424309)
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Matching Seq File Key = TRANSIT_0
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel stateID = key.charAt(8) = 0
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Transition Matrix (0, 0) = 7.212979465999484E-4
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Transition Matrix (0, 1) = 0.0029750348531350887
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Transition Matrix (0, 2) = 0.9963036672002649
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Matching Seq File Key = TRANSIT_1
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel stateID = key.charAt(8) = 1
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Transition Matrix (1, 0) = 0.011301229755584162
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Transition Matrix (1, 1) = 0.9886987696844521
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Transition Matrix (1, 2) = 5.59963649551057E-10
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Matching Seq File Key = TRANSIT_2
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel stateID = key.charAt(8) = 2
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Transition Matrix (2, 0) = 3.214352307450298E-10
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Transition Matrix (2, 1) = 0.5166010006877827
17/12/23 02:57:06 INFO BaumWelchUtils: CreateHmmModel Transition Matrix (2, 2) = 0.48339899899078204

```

Figure 6: End of of Baum Welch program execution

Table 2: List of the BaumWelch program options

No	Option	Argument(s)	Description	Custom Created
1	-input (-i)	input	Path to job input direc- tory.	No
2	-output (-o)	output	The directory path name for output.	No
3	- hmmModelInputPath (-hmmmin)	hmmModelInputPath [hmmModelInput- Path ...]	The input path of HMM Model. Must be of Se- quence Filetype.	Yes
4	-numHiddenStates (-numHid)	numHiddenStates [numHiddenStates ...]	Number of Hidden States of HMM	Yes
5	- numObservedStates (-numObs)	numObservedStates [numObserved- States ...]	Number of Observed States of HMM	Yes
6	-convergenceDelta (-cd)	convergenceDelta	The convergence delta value. Default is 0.5	No
7	-maxIter (-x)	maxIter	The maximum number of iterations.	No
8	-help (-h)		Print out help	No
9	-tempDir	tempDir	Intermediate output di- rectory	No
10	-startPhase	startPhase	First phase to run	No
11	-endPhase	endPhase	Last phase to run	No