

Hidden Markov Methods. Algorithms and Implementation

Final Project Report. MATH 127.

Nasser M. Abbasi

Course taken during Fall 2002 page compiled on July 2, 2015 at 12:08am

Contents

1	Example HMM	5
2	Forward α algorithm	7
2.1	Calculation of the alpha's using logarithms	9
3	Backward β algorithm	11
4	Viterbi algorithm	13
5	Configuration file description	15
6	Configuration file examples	16
6.1	Example 1. Casino example	16
6.2	Example 2. Homework problem 4.1	16
7	SHMM User interface	16
8	Source code description	17
9	Installation instructions	17

Abstract

Hidden Markov Models (HMM) main algorithms (forward, backward, and Viterbi) are outlined, and a GUI based implementation (in MATLAB) of a basic HMM is included along with a user guide. The implementation uses a flexible plain text configuration file format for describing the HMM.

Contents

Introduction

This report is part of my final term project for the course MATH 127, Mathematics department, Berkeley, CA, USA

I have taken this course as an extension student in the fall of 2002 to learn about computational methods in biology. The course is titled “Mathematical and Computational Methods in Molecular Biology” and given by Dr. Lior Pachter.

This report contains a detailed description of the three main HMM algorithms. An implementation (in MATLAB) was completed. See page 15 for detailed description of the implementation, including a user guide.

In the area of computational biology, HMMs are used in the following areas:

1. Pattern finding (Gene finding[3]).
2. multiple sequence alignments.
3. protein or RNA secondary structure prediction.
4. data mining and classification.

Definitions and notations

To explain HMM, a simple example is used.

Assume we have a sequence. This could be a sequence of letters, such as DNA or protein sequence, or a sequence of numerical values (a general form of a signal) or a sequence of symbols of any other type. This sequence is called the observation sequence.

Now, we examine the system that generated this sequence. Each symbol in the observation sequence is generated when the system was in some specific state. The system will continuously flip-flop between all the different states it can be in, and each time it makes a state change, a new symbol in the sequence is emitted.

In most instances, we are more interested in analyzing the system states sequence that generated the observation sequence. To be able to do this, we model the system as a stochastic model which when run, would have generated this sequence. One such model is the hidden Markov model. To put this in genomic perspective, if we are given a DNA sequence, we would be interested in knowing the structure of the sequence in terms of the location of the genes, the location of the splice sites, and the location of the exons and intron among others.

What we know about this system is what states it can be in, the probability of changing from one state to another, and the probability of generating a specific type of element from each state. Hence, associated with each symbol in the observation sequence is the state the system was in when it generated the symbol.

When looking a sequence of letter (a string) such as “ACGT”, we do not know which state the system was in when it generated, say the letter ‘C’. If the system have n states that it can be in, then the letter ‘C’ could have been generated from any one of those n states.

From the above brief description of the system, we see that there are two stochastic processes in progress at the same time. One stochastic process selects which state the system will switch to each time, and another stochastic process selects which type of element to be emitted when the system is in some specific state.

We can imagine a multi-faced coin being flipped to determine which state to switch to, and yet another such coin being flipped to determine which symbol to emit each time we are in a state.

In this report, the type of HMM that will be considered will be time-invariant. This means that the probabilities that describe the system do not change with time: When an HMM is in some state, say k , and there is a probability p to generate the symbol x from this state, then this probability do not change with time (or equivalent, the probability do not change with the length of the observation sequence).

In addition, we will only consider first-order HMM.¹

An HMM is formally defined² as a five-tuple:

1. Set of states $\{k, l, \dots, q\}$. These are the finite number of states the system can be in at any one time.

It is possible for the system to be in what is referred to as *silent states* in which no symbol is generated. Other than the special start and the end silent states (called the 0 states), all the system states will be active states (meaning the system will generate a symbol each time it is in one of these states).

When we want to refer to the state the system was in when it generated the element at position i in the observation sequence x , we write π_i . The initial state, which is the silent state the system was in initially, is written as π_0 . When we want to refer to the name of the state the system was in when it generated the symbol at position i we write $\pi_i = k$.

¹first-order HMM means that the current state of the system is dependent only on the previous state and not by any earlier states the system was in.

²In this report, I used the notations for describing the HMM algorithms as those used by Dr Lior Pachter, and not the notations used by the text book Durbin et al[1]

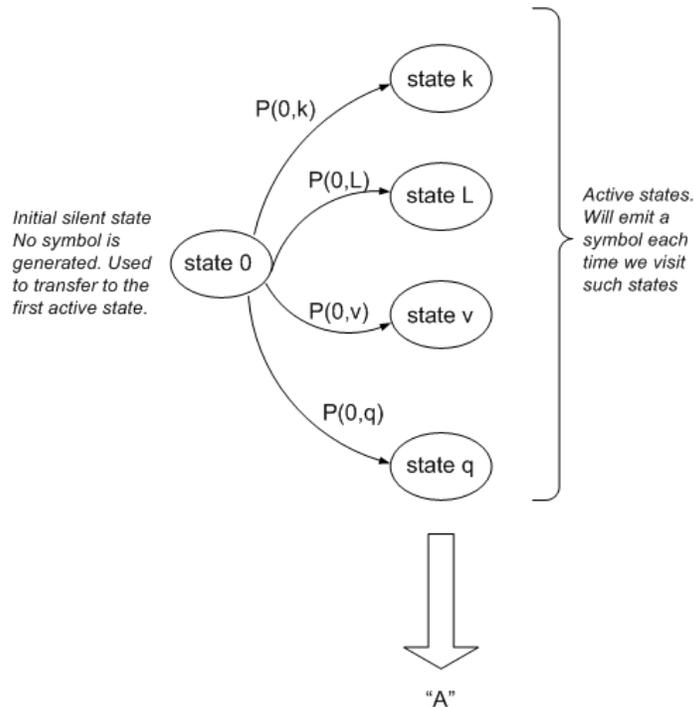


Figure 1: Initial silent state and initial transition to one of the active states. `silentState.eps`

2. P_{kl} is the probability of the system changing from state k to state l .

For a system with n number of states, there will be a transition matrix of size $n \times n$ which gives the state transition probability P_{kl} between any 2 states.

3. $b_k(i)$ is the emission probability. It is the probability of emitting the symbol seen at position i in the observation sequence x when the system was in state k .

As an example of the emission probability, assume that the observed sequence is $x = ATATTCGTC$ and assume that the system can be in any one of two states $\{k, l\}$. Then we write the probability of emitting the first symbol A when the system is in state k as $b_k(1)$.

We could also write the above as $b_k(A)$, since the symbol A is located at the first position of the sequence. In the above observation sequence, since the letter A occurred in the first and third positions, then $b_k(1)$ will have the same value as $b_k(3)$, since in both cases, the same type of letter is being emitted, which is the symbol A .

4. Set of initial probabilities P_{0k} for all states k . We imagine a special silent start state 0 from which the initial active state is selected according to the probability P_{0k} . see figure 1 on page 3 for illustration.

For a system with n number of states, there will be n such $P_{0,k}$ probabilities. So for some state, say q , the probability of this state being the first state the system starts in is given by P_{0q} . Notice that another notation for the initial probability P_{0k} is π_k , and this is used in [1] and other publication.

5. The set of possible observation sequences x . We imagine the system running and generating an observation sequence each time according to the emission probability distribution and according to the state sequence associated with this observation sequence. The collection of possible observation sequences is the language of the system. The unique alphabets that make up the observation sequences x are the allowed alphabet of the language. As an example, an HMM used to model DNA sequences would have $ACGT$ as its alphabets, while one used to model protein sequences would have the 20 amino acids as its alphabets.

So, to summarize how this model generates an observation sequence: The system starts in state 0. Then based the probability distribution P_{0k} called the initial state probability, a state k is selected. Then based on the

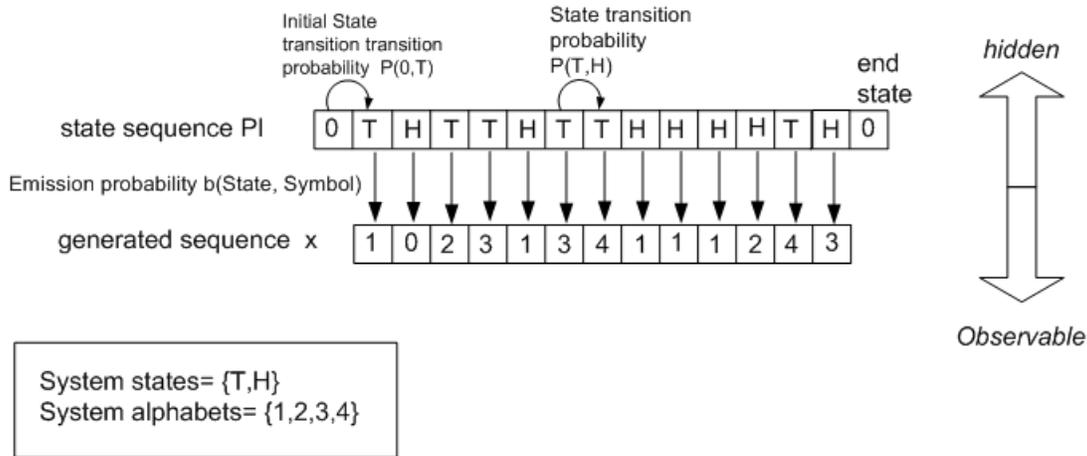


Figure 2: HMM parameters and how they are related. stateSeq.eps

emission probability distribution $b_k(i)$ for this state, the symbol at position i is selected and generated. Then a new state l is selected based on the transition probability distribution P_{kl} and the process repeats again as above. This process is continued until the last symbol is generated.

The sequence generated by the system is called the observation sequence x . The symbol at position i in the sequence is called $x[i]$ or x_i . We write $x[i] = A$ to mean that the symbol at position i is A .

In the implementation of HMM for this project, the length of each output symbol is limited to one character.

For a system with n states, and an observation sequence of length m , the number of possible state sequences that could have generated this observation sequence is n^m . Each one of these possible state sequences $\pi_{1\dots m}$ is of length m since the system will generate one symbol for each state transition as we are considering that there will not be any silent states other than the start and end states.

As an example, given a system with $n = 2$ with states $\{k, l\}$, and an observation sequence of ‘GCT’ of length $m = 3$, then there are $n^m = 2^3$ or 8 possible state sequences. They are

$$kkk, kkl, klk, kll, lkk, lkl, llk, lll$$

One of these 8 state sequences is the most likely one to have occurred.³

To summarize the above discussion, imagine we have a system with states $\{T, H\}$, and that the possible symbols or alphabets of the system it can generate are 1, 2, 3, 4. Figure 2 on page 4 shows one such possible state sequence π and the associated observation sequence x .

The questions one could ask about the above HMM system are:

1. Given a sequence of symbols x of length L that represents the output of the system (the observation sequence), determine the probability of the system generating this sequence up to and including $x[L]$. This is solved using the forward algorithm α .
2. Given a sequence of symbols x that represents the output of the system (the observations), determine the probability that some specific symbol in the observation sequence was generated when the system was in some specific state.

This is solved using the forward-backward algorithm ($\alpha\beta$). As an example, we can ask about the probability that some DNA nucleotide from a large DNA sequence was part of an exon, where an exon is one of the possible system states.

3. Given a sequence of symbols x of length L that represents the output of the system (the observation sequence), determine the most likely state transition sequence the system would have undergone in order to generate this observation sequence.

³One of the HMM algorithms, called Viterbi, is used to find which one of these sequences is the most likely one.

This is solved using the Viterbi algorithm γ . Note again that the sequence of states $\pi_{i\dots j}$ will have the same length as the observation sequence $x_{i\dots j}$.

The above problems are considered HMM analysis problems. There is another problem, which is the training or identification problem whose goal is to estimate the HMM parameters (the set of probabilities) we discussed on page 2 given some set of observation sequences taken from the type of data we wish to use the HMM to analyze. As an example, to train the HMM for gene finding for the mouse genome, we would train the HMM on a set of sequences taken from this genome.

One standard algorithm used for HMM parameter estimation (or HMM training) is called Baum-Welch, and is a specialized algorithm of the more general algorithm called EM (for expectation maximization). In the current MATLAB implementation, this algorithm is not implemented, but could be easily added later if time permits.

A small observation: In these notations, and in the algorithms described below, I use the value 1 as the first position (or index) of any sequence instead of the more common 0 used in languages such as C, C++ or Java. This is only for ease of implementation as I will be using MATLAB for the implementation of these algorithms, and MATLAB uses 1 and not 0 as the starting index of an array. This made translating the algorithms into MATLAB code much easier.

Graphical representations of an HMM

There are two main methods to graphically represent an HMM. One method shows the time axis (or the position in the observation sequence) and the symbols emitted at each time instance (called the trellis diagram). The other method is a more concise way of representing the HMM, but does not show which symbol is emitted at what time instance.

To describe each representation, I will use an example, and show how the example can be represented graphically in both methods.

1 Example HMM

Throughout this report, the following HMM example is used for illustration. The HMM configuration file for the example that can be used with the current implementation is shown. Each different HMM state description (HMM parameters) is written in a plain text configuration file and read by the HMM implementation. The format of the HMM configuration file is described on page 15.

Assume the observed sequence is *ATACC*. This means that *A* was the first symbol emitted, seen at time $t = 1$, and *T* is the second symbol seen at time $t = 2$, and so forth. Notice that we imagine that the system will write out each symbol to the end of the current observation sequence.

Next, Assume the system can be in any one of two internal states $\{S, T\}$.

Assume the state transition matrix P_{kl} to be as follows

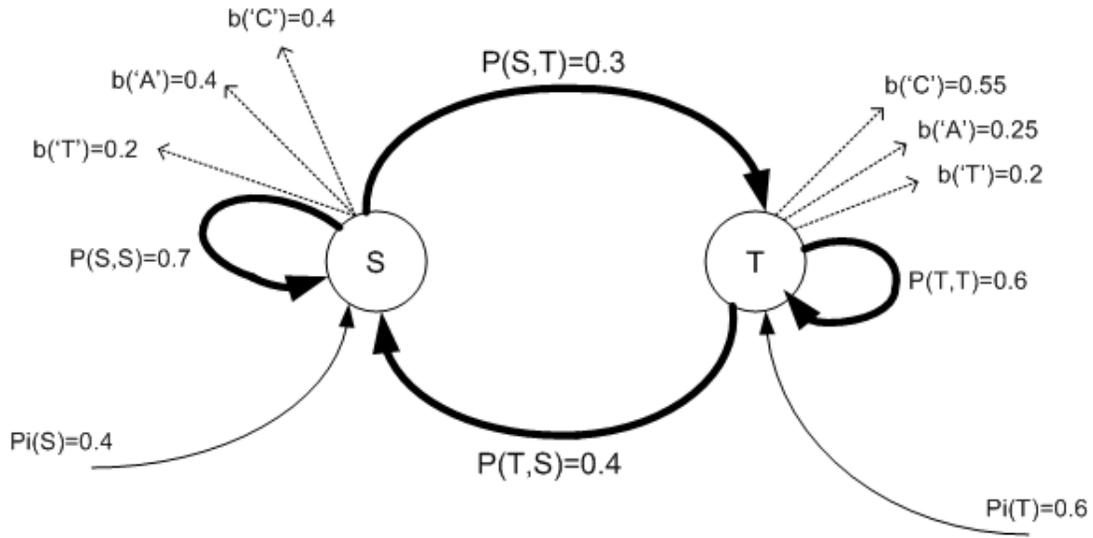
$$\begin{aligned} P_{ST} &= 0.3 & P_{SS} &= 0.7 \\ P_{TS} &= 0.4 & P_{TT} &= 0.6 \end{aligned}$$

And for the initial probabilities P_{0k} , assume the following

$$\begin{aligned} \pi_S &= 0.4 \\ \pi_T &= 0.6 \end{aligned}$$

And finally for the emission probabilities $b_k(i)$, there are 3 unique symbols that are generated in this example (the alphabets of the language of the observation sequence), and they are *A, T, C*. So there will be 6 emission probabilities, 3 symbols for each state, and given we have 2 states, then we will have a total of 6 emission probabilities. Assume they are given by the following emission matrix:

$$\begin{aligned} b_S(A) &= 0.4 & b_S(C) &= 0.4 & b_S(T) &= 0.2 \\ b_T(A) &= 0.25 & b_T(C) &= 0.55 & b_T(T) &= 0.2 \end{aligned}$$



HMM graphical description
of the given example.

Figure 3: HMM representation using node based diagram.

The first (non-time dependent) graphical representation of the above example is shown in figure 3 on page 6

In this representation, each state is shown as a circle (node). State transitions are represented by a directed edge between the states. Emission of a symbol is shown as an edge leaving the state with the symbol at the end of the edge.

Initial probability $\pi_{0,k}$ are shown as edges entering into the state from state 0 (not shown). We imagine that there is a silent start state 0 which all states originate from. The system cannot transition to state 0, it can only transition out of it.

The second graphical representation is a more useful way to represent an HMM for the description of the algorithms. This is called the trellis diagram.

Draw a grid, where time flows from left to right. On top of the grid show the symbols being emitted at each time instance. Each horizontal line in this grid represents one state, with state transitions shown as sloping edges between the nodes in the grid.

See figure 4 on page 7 for the graphical representation for the same example above using the trellis representation.

On each line segment in the grid, we draw an arrow from position i to $i + 1$ and label this arrow by the state transition probability from the state where the start of the arrow was to the state where the end of the arrow is. The weight of the edge is the state transition probability.

Notice that the edge can only be horizontal or sloped in direction. We can not have vertical edges, since we must advance by one time step to cause a state transition to occur. This means there will be as many internal states transitions as the number of time steps, or equivalently, as the number of symbols observed or the length of the observation sequence.

To use the above example in SHMM, we need to write down the state description. The format of writing down the state description is given on page 15. This below is the SHMM configuration file for the above example.

```
<states>
S
T
<init_prob>
0.4
0.6
```

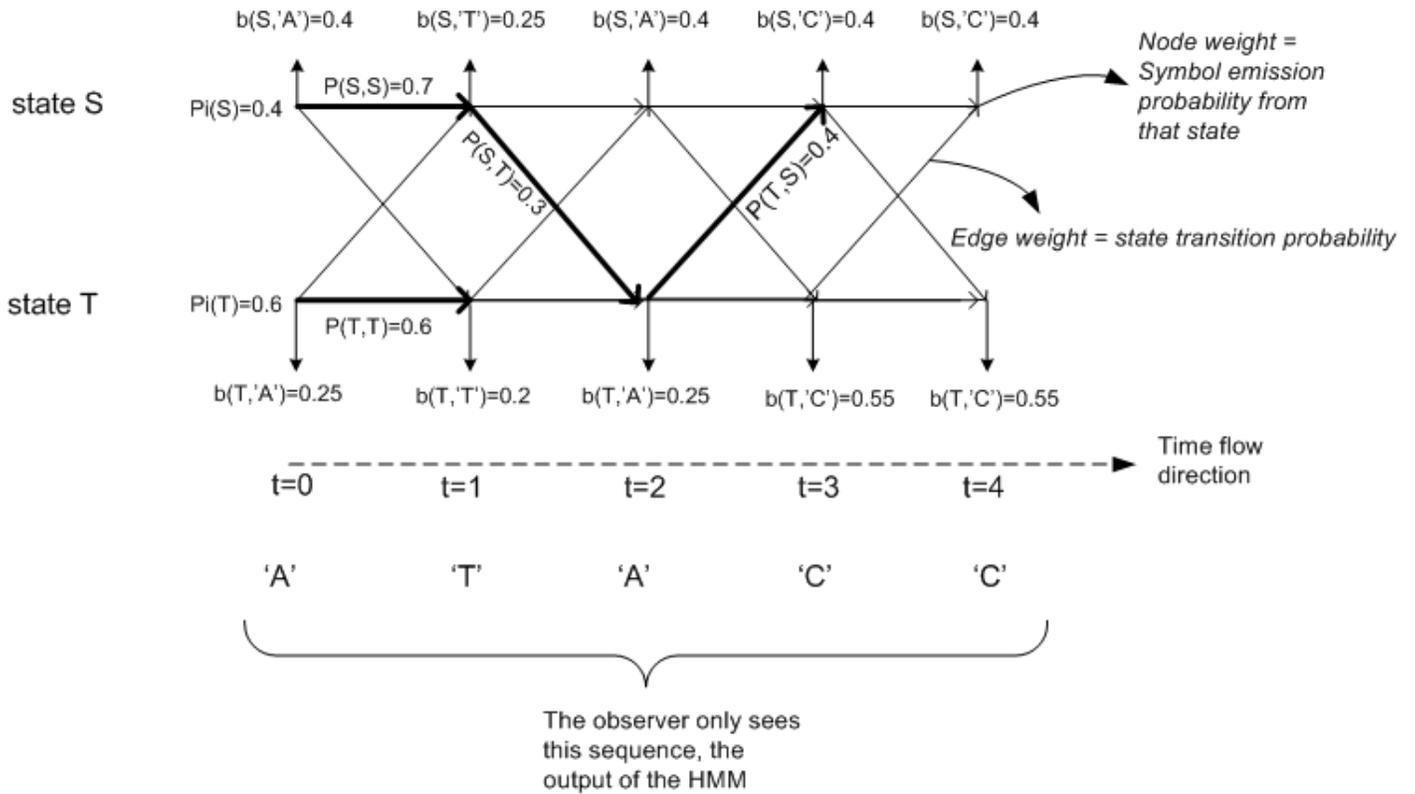


Figure 4: HMM representation using trellis diagram.

```

<symbols>
A,C,T
<emit_prob>
#emit probability from S state
0.4,0.4,0.2
#emit probability from T state
0.25,0.55,0.2

<tran_prob>
0.7,  0.3
0.4,  0.6

```

HMM algorithms

Given the above notations and definitions, we are ready to describe the HMM algorithms.

2 Forward α algorithm

Here we wish to find the probability p that some observation sequence x of length m was generated by the system which has n states. We want to find the probability of the sequence up to and including $x[m]$, in other words, the probability of the sequence $x[1] \cdots x[m]$.

The most direct and obvious method to find this probability, is to enumerate all the possible state transition sequences $q = \pi_{1..m}$ the system could have gone through. Recall that for an observation sequence of length m there will be n^m possible state transition sequences.

For each one of these state transition sequences q_i (that is, we fix the state transition sequence), we then find the probability of generating the observation sequence $x[1] \cdots x[m]$ for this one state transition sequence.

After finding the probability of generating the observation sequence for each possible state transition sequence,

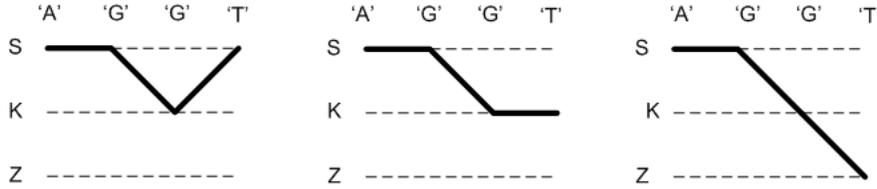


Figure 5: Finding the probability of a an observation sequence using brute search method.

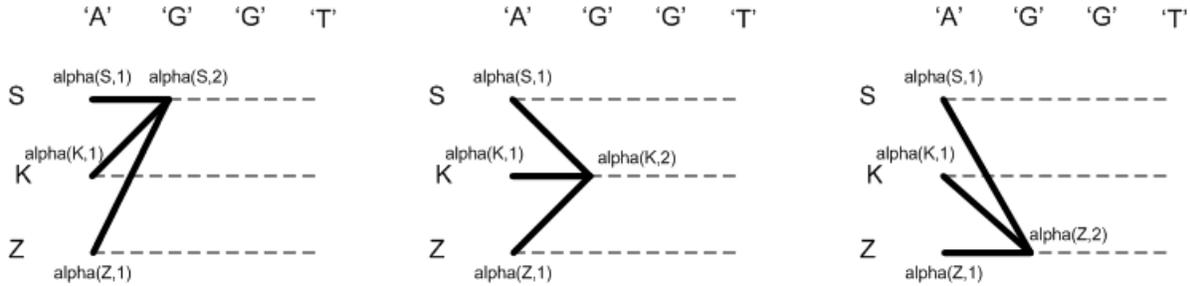


Figure 6: Finding the probability of a sequence using the α method.

we sum all the resulting probabilities to obtain the probability of the observation sequence being generated by the system.

As an example, for a system with states $\{S, K, Z\}$ and observation sequence $AGGT$ a partial list of the possible state sequences we need to consider is shown in figure 5 on page 8.

We can perform the above algorithm as follows:

Find the probability of the system emitting the symbol $x[1]$ from π_1 . This is equal to the probability of emitting the symbol $x[1]$ from state π_1 multiplied by the probability of being in state π_1 .

Next, find the probability of generating the sequence $x[1] \cdots x[2]$. This is equal to the probability of emitting the sequence $x[1]$ (which we found from above), multiplied by the probability of moving from state π_1 to state π_2 , multiplied by the probability of emitting the symbol $x[2]$ from state π_2 .

Next, find the probability of generating the sequence $x[1] \cdots x[3]$. This is equal to the probability of emitting the sequence $x[1] \cdots x[2]$ (which we found from above), multiplied by the probability of moving from state π_2 to state π_3 , multiplied by the probability of emitting the symbol $x[3]$ from state π_3 .

We continue until we find the probability of emitting the full sequence $x[1] \cdots x[m]$

Obviously the above algorithm is not efficient. There are n^m possible state sequences, and for each one we have to perform $O(2m)$ multiplications, resulting in $O(m n^m)$ multiplications in total.

For an observation sequence of length $m = 500$ and a system with $n = 10$ states, we have to wait a very long time for this computation to complete.

A much improved method is the α algorithm. The idea of the α algorithm, is that instead of finding the probability for each different state sequence separately and then summing the resulting probabilities, we scan the sequence $x[1] \cdots x[i]$ from left to right, one position at a time, but for each position we calculate the probability of generating the sequence up to this position for all the possible states.

Let us call the probability of emitting the sequence $x[1] \cdots x[i]$ as $\alpha(i)$. And we call the probability of emitting $x[1] \cdots x[i]$ when the system was in state k when emitting $x[i]$ as $\alpha_k(i)$.

So, the probability of emitting the sequence $x[1] \cdots x[i]$ will be equal to the probability of emitting the same sequence but with $x[i]$ being emitted from one state say k , this is called $\alpha_k(i)$, plus the probability of emitting the same sequence but with $x[i]$ being emitted from another state, say l , this is called $\alpha_l(i)$, and so forth until we have covered all the possible states. In other words $\alpha(i)$ is the $\sum_k \alpha_k(i)$ for all states k .

Each $\alpha_k(i)$ is found by multiplying the $\alpha(i-1)$ for each of the states, by the probability of moving from each of those states to state k , multiplied by the probability of emitting $x[i]$ from state k . This step is illustrated in figure 6 on page 8.

Each time we calculate $\alpha_k(i)$ for some state k , we store this value, and use it when calculating $\alpha(i+1)$.

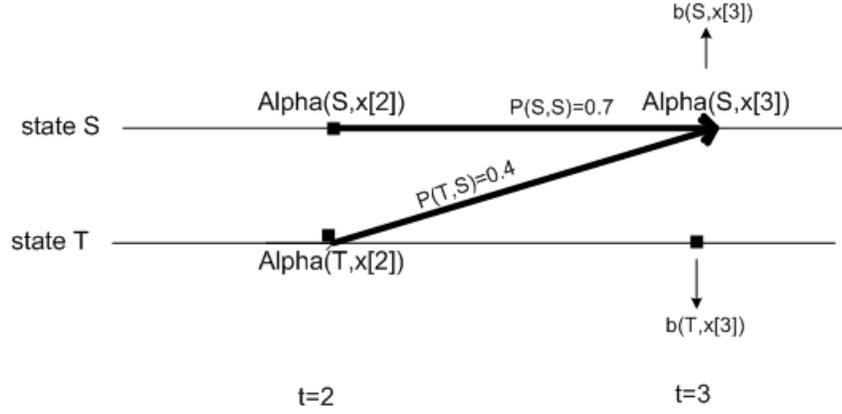


Figure 7: Recursive step for the α method for state S.

Hence We only need to store the α 's for one position earlier since we are working with a first-order HMM here. This means we need space requirements for two columns each of length of the number of states n .

From above, we see that we have reduced the running time $O(mn^m)$ of the direct method, to $O(mn)$ with the α method.

The following is a more formal description of the algorithm to find α for one state, say state k .

Initialization step:

$$\alpha_k(1) = p_{0k}b_k(1)$$

What the above says is that α , initially (at position $i = 1$) and for state k , is equal to the probability of the HMM starting initially in state k multiplied by the probability of emitting the symbol found at $x[1]$ from state k .

Iterative step:

$$j = 2 \cdots i \quad \alpha_k(j) = b_k(j) \sum_l \alpha_l(j-1)P_{lk}$$

What the above says, is that α at position $i = j$ when in state k , is the sum of all the α 's from one observation earlier collected from all the states, multiplied by the probability of transition from each of those states to the state k , multiplied by the probability of emitting the symbol x_j when in state k .

I will illustrate this iterative step for one time step, using the same example from above, and using the grid representation of an HMM to find $\alpha(3)$. Since there are 2 states $\{S, T\}$ in this example, we need to calculate the iterative step for each state.

When the system is in state S at $j = 3$ we obtain the diagram shown in figure 7 on page 9.

In the above example, to find $\alpha_S(3)$, we obtain

$$\alpha_S(3) = b_S(3)[\alpha_S(2)P_{SS} + \alpha_T(2)P_{TS}]$$

$\alpha_T(3)$ is found similarly, and is shown in figure 8 on page 10.

In the above example, for $\alpha_T(3)$ we obtain

$$\alpha_T(3) = b_T(3)[\alpha_T(2)P_{TT} + \alpha_S(2)P_{ST}]$$

Hence, to find the probability of the sequence up to and including x_3 , the final answer would be

$$\begin{aligned} \alpha(x_3) &= \alpha_T(3) + \alpha_S(3) \\ &= b_S(3)[\alpha_S(2)P_{SS} + \alpha_T(2)P_{TS}] + b_T(3)[\alpha_T(2)P_{TT} + \alpha_S(2)P_{ST}] \end{aligned}$$

2.1 Calculation of the alpha's using logarithms

Due to numerical issues (underflow) which will occur when multiplying probabilities (values less than 1.0) repeatedly as we calculate the α in the forward direction, the log of the probabilities will be used instead, and

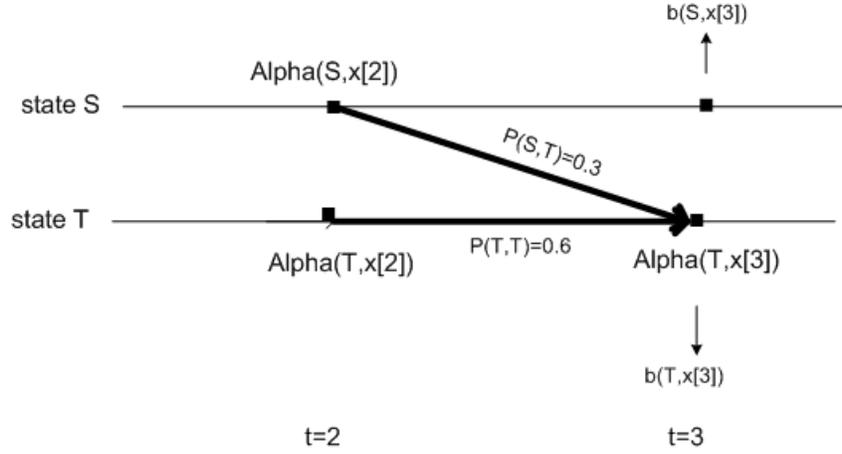


Figure 8: Recursive step for the α method for state T.

multiplications become summations. Any logarithm for a base greater than 1 can be used. In the MATLAB implementation, I used the natural logarithms e .

Initially, all the probabilities P (this includes the initial state probabilities $\pi_{(0,k)}$, the state transition probabilities matrix P_{xy} , and the emission probabilities matrix $b_S(i)$) are converted to $\log(P)$ and stored in cache for later access, this speeds the implementation by not having to call the log function each time the program is run for different observation sequence but for the same HMM.

Assume the system has states $\{A, B, C, \dots, Z\}$, and we wish to find the α for the observation sequence up to and including $x[i]$ when the state at i was T . We obtain

$$\alpha_T(i) = b_T(i)[\alpha_A(i-1)P_{AT} + \alpha_B(i-1)P_{BT} + \dots + \alpha_Z(i-1)P_{ZT}]$$

Taking the log of the above results in

$$\log(\alpha_T(i)) = \log(b_T(i)) + \log[\alpha_A(i-1)P_{AT} + \alpha_B(i-1)P_{BT} + \dots + \alpha_Z(i-1)P_{ZT}]$$

Now, looking at the term $\log[\alpha_A(i-1)P_{AT} + \dots + \alpha_Z(i-1)P_{ZT}]$, this is the same as $\log(x + y + \dots + z)$. Here, we know the values of $\log(x), \log(y), \dots, \log(z)$, and not x, y, \dots, z (recall that we have converted all the probabilities to log initially for performance purposes). So we need to find $\log(x + y)$ given $\log(x)$ and $\log(y)$.

This is done using the following transformation

$$\begin{aligned} \log(x + y) &= \log\left(x\left(1 + \frac{y}{x}\right)\right) \\ &= \log\left(x\left(1 + e^{\log\left(\frac{y}{x}\right)}\right)\right) \\ &= \log(x) + \log\left(1 + e^{\log(y) - \log(x)}\right) \end{aligned}$$

Applying the above transformation to many terms being summed, and not just two, is straightforward. Take x as the first term, and let y be the rest of the terms (the tail of the expression), we then apply the above transformation recursively to the tail to compute the log of the sum of the probabilities as long as the tail is of length greater than one element. As shown below.

$$\log(a + b + c + \dots + y + z) = \log(a + H) = \log(a) + \log\left(1 + e^{\log(H) - \log(a)}\right)$$

Where

$$H = (b + c + \dots + y + z)$$

Now repeat the above process to find $\log(b + c + \dots + y + z)$.

$$\log(b + c + \dots + y + z) = \log(b + H) = \log(b) + \log(1 + e^{\log(H) - \log(b)})$$

Where now $H = (c + \dots + y + z)$.

Continue the above recursive process, until the expression H contains one term to obtain

$$\log(y + z) = \log(y) + \log(1 + e^{\log(z) - \log(y)})$$

This is a sketch of a function which accomplishes the above.

```
function Log_Of_Sums(A: array 1..n of probabilities) returns double IS
begin
  IF (length(A)>2) THEN
    return( log(A[1]) + log( 1 + exp( Log_of_Sums(A[2..n])) - log( A[1] ) ));
  ELSE
    return( log(A[1] + log( 1 + exp( A[2] - log( A[1] ) )));
  END
end log_Of_Sums
```

The number of states is usually much less than the number of observations, so there is no problems (such as stack overflow issues) with using recursion in this case.

3 Backward β algorithm

Here we are interested in finding the probability that some part of the sequence was generated when the system was in some specific state.

For example, we might want to find, looking at a DNA sequence x , the probability that some nucleotide at some position, say $x[i]$ was generated from an exon, where an exon is one of the possible states the system can be in.

Let such state be S . Let the length of the sequence be L , then we need to obtain $P(\pi_i = S | x_{i..L})$. In the following, when x is written without subscripts, it is understood to stand for the whole sequence $x_1 \dots x_L$. Also x_i is the same as $x[i]$.

Using Baye's probability rule we get

$$P(\pi_i = S | x) = \frac{P(\pi_i = S, x)}{P(x)}$$

But the joint probability $P(\pi_i = S, x)$ can be found from

$$P(\pi_i = S, x) = P(x_1 \dots x_i, \pi_i = S) P(x_{i+1} \dots x_L | \pi_i = S)$$

Looking at the equation above, we see that $P(x_1 \dots x_i, \pi_i = S)$ is nothing but $\alpha_S(i)$, which is defined as the probability of the sequence $x_1 \dots x_i$ with $x[i]$ being generated from state S .

The second quantity $P(x_{i+1} \dots x_L | \pi_i = S)$ is found using the backward β algorithm as described below.

Initialization step:

For each state k , let $\beta_k(L) = 1$

recursive step:

$$i = (L - 1) \dots 1 \quad \beta_k(i) = \sum_l P_{kl} b_l(i + 1) \beta_l(i + 1)$$

Where the sum above is done over all states. Hence

$$P(\pi_i = S | x) = \frac{\alpha_S(i) \beta_S(i)}{P(x)}$$

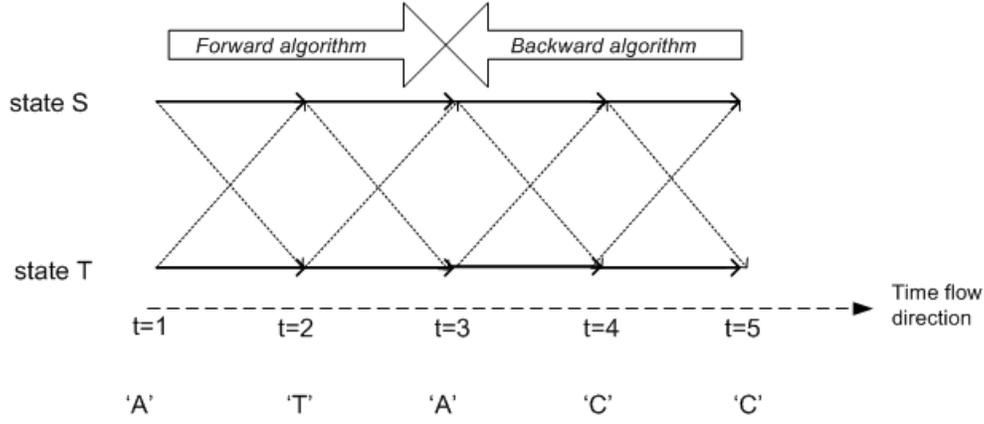


Figure 9: The forward-backward algorithm.

Where $P(x)$ is found by running the forward algorithm to the end of the observed sequence x , so $P(x) = \alpha(L)$. Hence we obtain

$$P(\pi_i = S|x) = \frac{\alpha_S(i)\beta_S(i)}{\alpha(L)}$$

To illustrate this algorithm, Using the grid HMM diagram, I will use the example on page 5, to answer the question of what is the probability of the observation at position 3 (which is the letter 'A') coming from state S, given the observation ATACC.

In the above, we will run the forward α algorithm from position 1 to position 3, then run the backward algorithm from the end of the sequence back to position 3. (In practice, the forward α 's and backward β 's are calculated once for all the states and all the positions (that is, for all the nodes in the trellis diagram) and stored for later lookup). See figure 9 on page 12 for an illustration.

Initialization step:

$$\begin{aligned}\beta_S(5) &= 1 \\ \beta_T(5) &= 1\end{aligned}$$

recursive step: $t = (5 - 1) \dots 1$.

At time $t = 4$, we get

$$\begin{aligned}\beta_S(4) &= P_{ST}b_T(5)\beta_T(5) + P_{SS}b_S(5)\beta_S(5) \\ \beta_T(4) &= P_{TT}b_T(5)\beta_T(5) + P_{TS}b_S(5)\beta_S(5)\end{aligned}$$

The above quantities are calculated using the parameters in this example, resulting in

$$\begin{aligned}\beta_S(4) &= 0.3 \times 0.55 \times 1 + 0.7 \times 0.4 \times 1 = 0.445 \\ \beta_T(4) &= 0.6 \times 0.55 \times 1 + 0.4 \times 0.4 \times 1 = 0.49\end{aligned}$$

at time $t = 3$ we get

$$\begin{aligned}\beta_S(3) &= P_{ST}b_T(4)\beta_T(4) + P_{SS}b_S(4)\beta_S(4) \\ \beta_T(3) &= P_{TT}b_T(4)\beta_T(4) + P_{TS}b_S(4)\beta_S(4)\end{aligned}$$

The above quantities are calculated using the parameters in this example, resulting in

$$\begin{aligned}\beta_S(3) &= 0.3 \times 0.55 \times 0.49 + 0.7 \times 0.4 \times 0.445 = 0.2054 \\ \beta_T(3) &= 0.6 \times 0.55 \times 0.49 + 0.4 \times 0.4 \times 0.445 = 0.2329\end{aligned}$$

And so forth, until we reach $t = 1$

$$\begin{aligned}\beta_S(1) &= P_{ST}b_T(2)\beta_T(2) + P_{SS}b_S(2)\beta_S(2) \\ \beta_T(1) &= P_{TT}b_T(2)\beta_T(2) + P_{TS}b_S(2)\beta_S(2)\end{aligned}$$

This completes the iterative step of the backward algorithm. From the above we see that $\beta_S(3) = 0.205449$ Now, apply the forward algorithm from position 1 to position 5, and find $\alpha_S(3)$.

I will not show this step here, since the forward algorithm is already described above. But if we calculate it, we will get $\alpha(S, t = 3) = 0.012488$

So

$$P(\pi_3 = S|ATACC) = \frac{\alpha_S(3)\beta_S(3)}{P(x)} = \frac{0.012488 \times 0.205449}{P(x)} = \frac{0.0025656}{P(x)}$$

Where $P(x)$ is the same as $\alpha(L)$, that is, the forward algorithm found at the end of the sequence x . $\alpha(L)$ can be found to be 0.0046025920, so

$$P(\pi_3 = S|ATACC) = \frac{0.0025656}{0.0046025920} = 0.557438$$

4 Viterbi algorithm

The Viterbi algorithm is used to find the most likely state transition sequence $q = \pi_1 \cdots \pi_i$ associated with some observation sequence $x = x_1 \cdots x_i$.

Since one symbol is generated each time the system switches to a new state, the state sequence q will be of the same length as that of the observation sequence x .

In this algorithm, we need to calculate γ for each position, where $\gamma_k(i)$ means the probability of the most likely state path ending in state k when generating the symbol at position i in the observation sequence.

If we know $\gamma_k(i)$ for each state k , then we can easily find $\gamma_l(i + 1)$ for each state l , since we just need to find the most probable state transition from each state k to each state l .

The following algorithm finds $\gamma_k(i)$ at position i for state k .

Initialization step:

Find γ at position $i = 1$ for all states. Assume we have n states, a, b, c, \dots, z , then

$$\begin{aligned}\gamma_a(1) &= P_{0a} b_a(1) \\ \gamma_b(1) &= P_{0b} b_b(1) \\ &\vdots \\ \gamma_z(1) &= P_{0z} b_z(1)\end{aligned}$$

Now we need to calculate γ for the rest of the sequence.

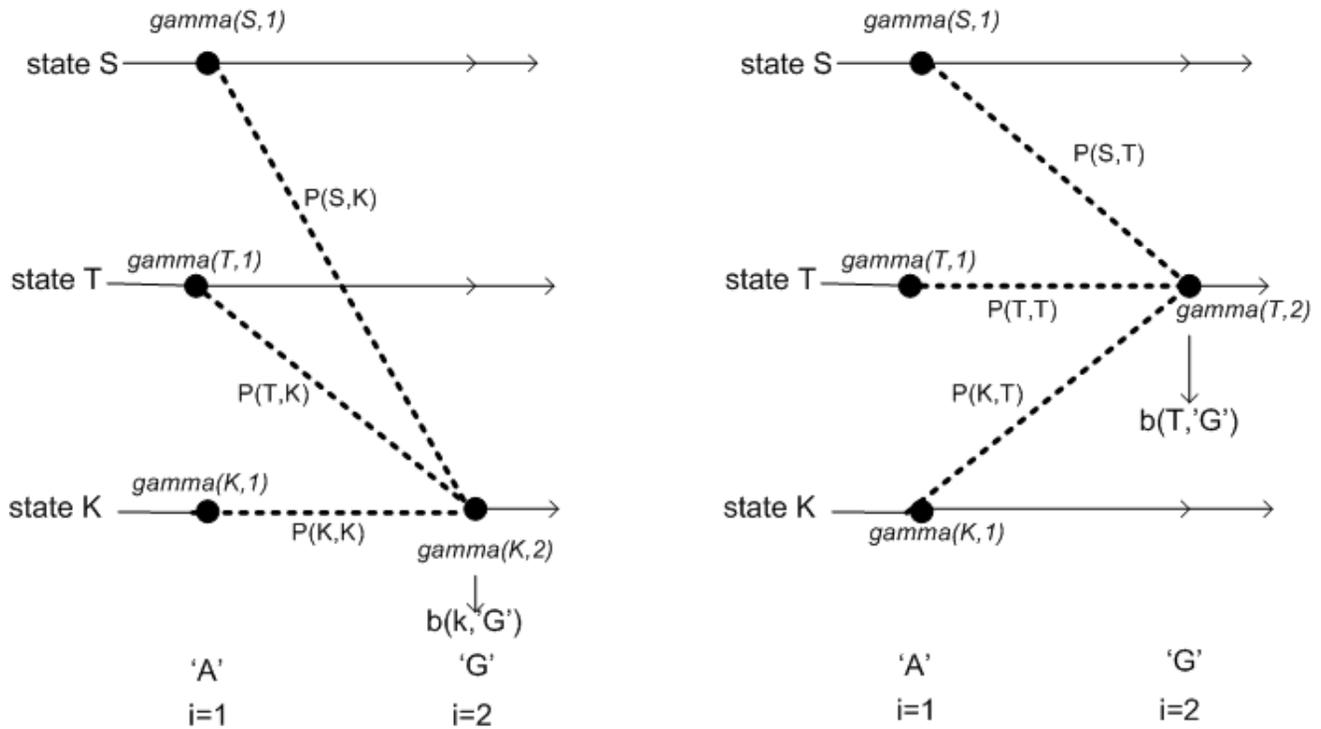
Assume we have states a, b, m, c, \dots, z , then to find γ for each of these states at each time step we do the following: (this below shows finding γ for some state m)

iterative step: $j = 2 \dots i$

$$\gamma_m(j) = b_m(j) \max \begin{cases} P_{am} \gamma_a(j-1) \\ P_{bm} \gamma_b(j-1) \\ P_{mm} \gamma_m(j-1) \\ \vdots \\ P_{zm} \gamma_z(j-1) \end{cases}$$

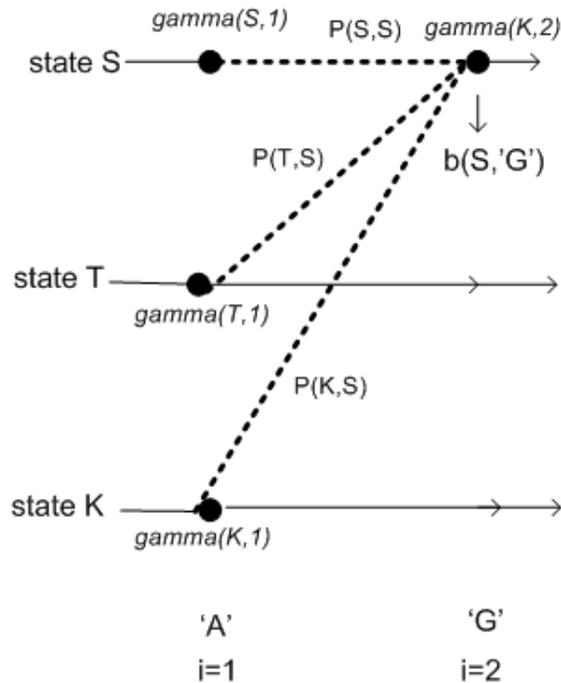
The above is illustrated in figure 10 on page 14.

Now that we have found γ for each state at each position $1 \cdots i$, we can find the most likely state sequence up to any position $j \leq i$ by finding which state had the largest γ at each position.



$$\begin{aligned} \text{Gamma}(K,2) &= b(k, 'G') * \\ &\text{MAX}(\\ &\quad \gamma(S,1) * P(S,K), \\ &\quad \gamma(T,1) * P(T,K), \\ &\quad \gamma(K,1) * P(K,K) \\ &) \end{aligned}$$

$$\begin{aligned} \text{Gamma}(T,2) &= b(T, 'G') * \\ &\text{MAX}(\\ &\quad \gamma(S,1) * P(S,T), \\ &\quad \gamma(T,1) * P(T,T), \\ &\quad \gamma(K,1) * P(K,T) \\ &) \end{aligned}$$



$$\begin{aligned} \text{Gamma}(S,2) &= b(S, 'G') * \\ &\text{MAX}(\\ &\quad \gamma(S,1) * P(S,S), \\ &\quad \gamma(T,1) * P(T,S), \\ &\quad \gamma(K,1) * P(K,S) \\ &) \end{aligned}$$

Figure 10: The Viterbi iterative step to find γ for each position viterbi.eps

Design and user guide of SHMM

To help understand HMMs more, I have implemented, in MATLAB, a basic HMM with a GUI user interface for ease of use. This tool accepts an HMM description (model parameters) from a plain text configuration file. The tool (called SHMM, for simple HMM) implements the forward, backward, and Viterbi algorithms. This implementation is for a first order HMM model.

Due to lack of time to implement an HMM for an order greater than a first order, which is needed to enable this tool to be used for such problems as gene finding (to detect splice sites, start of exon, etc. . .). This implementation therefor is restricted to one character look-ahead on the observation sequence. However, even with this limitation, this tool can be useful for learning how HMM's work.

SHMM takes as input a plain text configuration file (default .hmm extension), which contains the description of the HMM (states, initial probabilities, emission probabilities, symbols and transition probabilities). And the user is asked to type in the observation sequence in a separate window for analysis.

5 Configuration file description

The configuration file is broken into 5 sections. The order of the sections must be as shown below. All 5 sections must exist. Each section starts with the name of the section, contained inside a left < and a right > brackets. No space is allowed after < or before >. The entries in the file are case sensitive. So a state called 'sunny' will be taken as different from a state called 'SUNNY'. The same applies to the characters allowed in the observation sequence: the letter 'g' is taken as different from 'G'.

The sections, and the order in which they appear in the configuration file, must be as this, from top to bottom:

1. <states>
2. <init_prob>
3. <symbols>
4. <emit_prob>
5. <tran_prob>

In the <states> sections, the names of the states are listed. One name per line. Names must not contain a space.

In the <init_prob> section, the initial probabilities of being in each state are listed. One per line. One number per line. The order of the states is assumed the same as the order in the <states> section. This means the initial probability of being in the first state listed in the <states> section will be assumed to be the first number in the <init_prob> section.

In the <symbols> section, the observation symbols are listed. Each symbol must be one character long, separated by a comma. Can use as many lines as needed to list all the symbols. The symbol must be an alphanumeric character.

In the <emit_prob> section, we list the probability of emitting each symbol from each state. The order of the states is assumed to be the same order of the states as in the <states> section.

The emission probability for each symbol is read. The order of the symbols is taken as the same order shown in the symbols section. The emission probabilities are separated by a comma. One line per state. See examples below.

In the <tran_prob>, we list the transition probability matrix. This is the probability of changing state from each state to another. One line represent one row in the matrix. Entries are separated by a comma. One row must occupy one line only. The order of the rows is the same as the order of the states shown in the <states> section.

For example, if we have 2 states S1,S2, then the first line will be the row that represents S1->S1, and S1->S2 transition probabilities. Where S1->S1 is in position 1,1 of matrix, and S1->S2 will be in position 1,2 in the matrix. See example below for illustration.

Comments in the configuration file are supported. A comment can start with the symbol #. Any line that starts with a # is ignored and not parsed. The comment symbol # must be in the first position of the line.

To help illustrate the format of the HMM configuration file, I will show an annotated configuration file for the casino example shown in the book 'Biological sequence analysis' by Durbin et. al [1], page 54. In this example, we have 2 dies, one is a fair die, and one is a loaded die. The following is the HMM configuration file.

```
# example HMM configuration file for Durbin, page 54
# die example.

<states>
fair
loaded

<init_prob>
#fair init prob
0.5

#loaded init prob
0.5

<symbols>
1,2,3,4,5,6

<emit_prob>
#emit probability from fair state
1/6, 1/6, 1/6, 1/6, 1/6, 1/6
#emit probability from fair state
1/10, 1/10, 1/10, 1/10, 1/10, 1/10

<tran_prob>
0.95, 0.05
0.1, 0.9
```

See figure 11 on page 18 for description of each field in the HMM configuration file.

6 Configuration file examples

To help better understand the format of the HMM configuration file, I will show below a number of HMM states and with the configuration file for each.

6.1 Example 1. Casino example

This example is taken from page 54, Durbin et. al book [1]. See figure 12 on page 19

6.2 Example 2. Homework problem 4.1

This example is taken from Math 127 course, homework problem 4.1, Fall 2002.

See figure 13 on page 20

7 SHMM User interface

Figure 14 on page 21 shows the SHMM GUI. To start using SHMM, the first thing to do is to load the HMM configuration file which contains the HMM parameters. SHMM will parse the file and will display those

parameters. Next, type into the observation sequence window the observation sequence itself. This sequence could be cut and pasted into the window as well. Spaces and new lines are ignored.

Next, hit the top right RUN button. This will cause SHMM to evaluate the α , β and γ for each node in the trellis diagram, and will display all these values. For the γ calculations, the nodes in the trellis diagram which are on the Viterbi path have a '*' next to them.

To invoke the (α, β) algorithm, enter the position in the observation and the state name and hit the small RUN button inside the window titled **Alpha Beta**.

To modify the HMM parameters and run SHMM again, simply use an outside editor to do the changes in the HMM configuration file, then load the file again into SHMM using the 'Open' button.

8 Source code description

The following is a list of all the MATLAB files that make up SHMM, with a brief description of each file. At the end of this report is the source code listing.

- `nma_hmmGUI.fig`: The binary description of the GUI itself. It is automatically generated by MATLAB using GUIDE.
- `nma_hmmGUI.m`: The MATLAB code that contains the properties of each component in the GUI.
- `nma_hmm_main.m`: This is the main function, which loads the GUI and initializes SHMM.
- `nma_hmm_callbacks.m`: The callbacks which are called when the user activate any control on the GUI.
- `nma_readHMM.m`:
- `nma_hmmBackward.m`: The function which calculates β .
- `nma_hmmForward.m`: The function which calculates α .
- `nma_hmm_veterbi.m`: The function which calculates γ .
- `nma_hmmEmitProb.m`: A utility function which returns the emission probability given a symbol and a state.
- `nma_trim.m`: A utility function which removes leading and trailing spaces from a string.
- `getPosOfState.m`: A utility function which returns position of state in an array given the state name.
- `nma_logOfSums.m`: A utility function which calculates the log of sum of logs for the α computation.
- `nma_logOfSumsBeta.m`: A utility function which calculates the log of sum of logs for the β computation.

9 Installation instructions

Included on the CDROM are 2 files: `setup.exe` and `unix.tar`. The file `setup.exe` is the installer for windows. On windows, simply double click the file `setup.exe` and follow the instructions. When installation is completed, an icon will be created on the desktop called *MATH127*. Double click on this icon to open. Then 2 icons will appear. One of them is titled HMM. Double click on that to start SHMM. A GUI will appear.

On Unix (Linux or Solaris), no installer is supplied. Simply untar the file `math127.zip`. This will create the same tree that the installer would create on windows. The tree contains a folder called *src*, which contains the source code needed to run SHMM from inside matlab. To run SHMM on Unix, `cd` to the folder where you untarred the file `math127.tar`, which will be called *MATH127* by default, and navigate to the `src` folder for the HMM project. Then run the file `nma_hmm_main` from inside MATLAB console to start SHMM.

This report is in .ps format, and the original .tex file and the diagram files are included in the *doc* folder. Example HMM configuration files for use with SHMM are included in the *conf* folder.

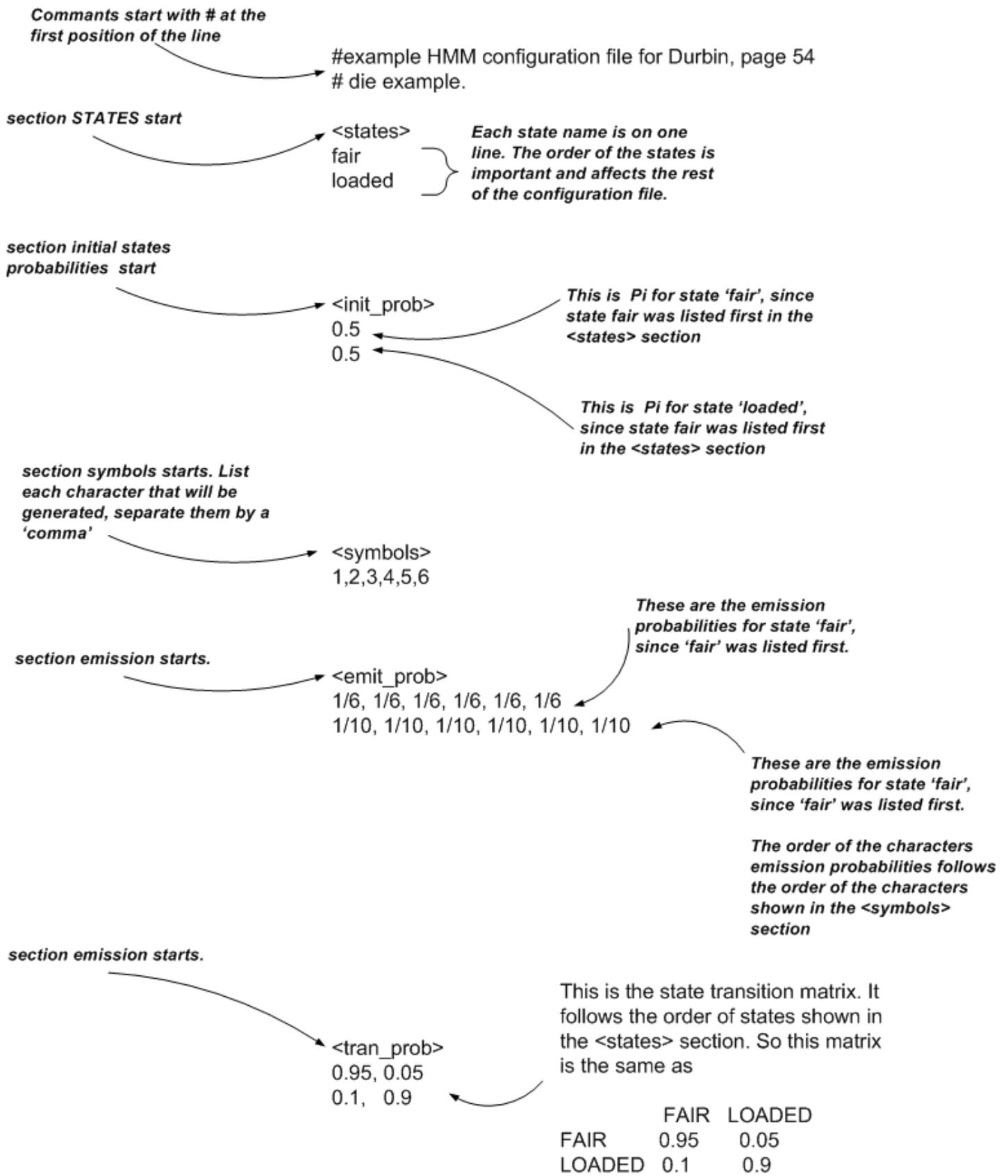
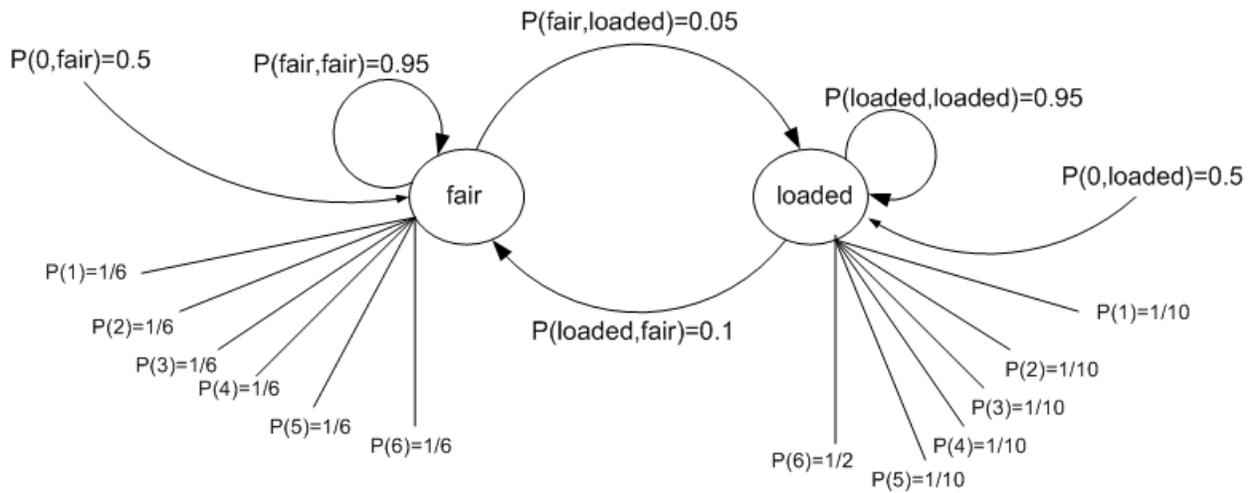


Figure 11: Description of fields in the HMM configuration file. The HMM parameters



```

<states>
fair
loaded

<init_prob>
#fair init prob
0.5

#loaded init prob
0.5

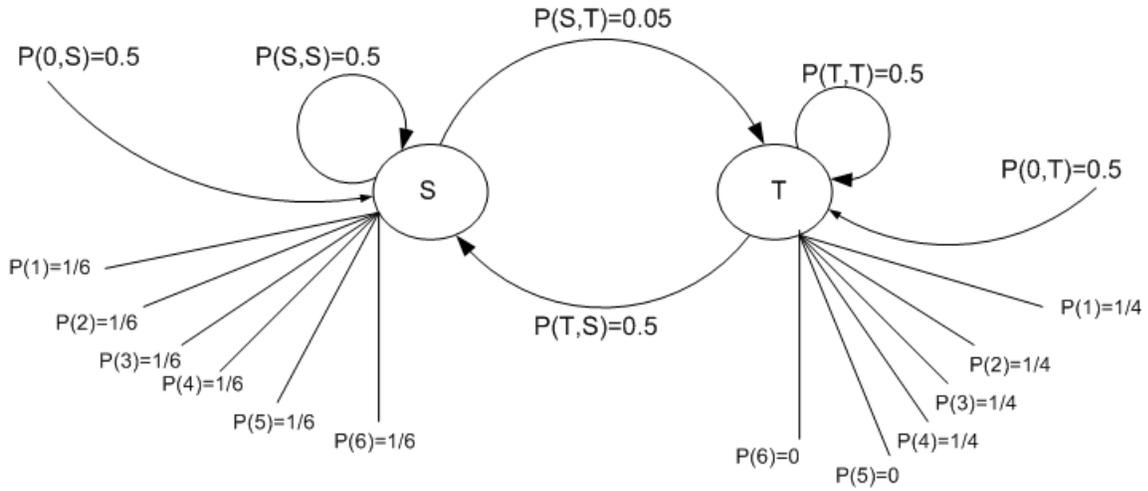
<symbols>
1,2,3,4,5,6

<emit_prob>
#emit probability from fair state
1/6, 1/6, 1/6, 1/6, 1/6, 1/6
#emit probability from fair state
1/10, 1/10, 1/10, 1/10, 1/10, 1/10

<tran_prob>
0.95, 0.05
0.1, 0.9

```

Figure 12: Example 1 HMM configuration file



```

<states>
S
T

<init_prob>
0.5
0.5

<symbols>
1, 2, 3,4,5,6

<emit_prob>
0.166666667, 0.166666667, 0.166666667, 0.166666667, 0.166666667, 0.166666667
0.25, 0.25, 0.25 , 0.25 , 0 , 0

<tran_prob>
0.5, 0.5
0.5, 0.5

```

Figure 13: Example 2 HMM configuration file

Click here to open HMM configuration file.

Type in, or cut/paste, the observation sequence here. (Spaces are ignored)

Click here to calculate the Alpha, Beta, and Gamma probabilities.

These are read and parsed from the HMM configuration file. They are the HMM model parameters. To change them, change the file and open it again.

These 2 input windows are used for the Alpha-Beta (Forward-Backward) calculation. Click on the small Run button to evaluate the Alpha-Beta probability.

The position line shows the position in the observation sequence.

This window display the Beta (backward) probabilities at each node in the trellis diagram.

This window display the Gamma (Viterbi) probabilities at each node in the trellis diagram.

This value is the probability of the observation sequence.

This window display the ALpha (Forward) probabilities at each node in the trellis diagram.

This window shows the Viterbi state path sequence.

A "*" here indicates that this state is on the most likely state path

Simple HMM model, V 1.0. Written by Nasser Abbasi for MATH 127 course taught by Dr Lior Pachter, UC Berkeley, Fall 2002.

Open HMM configuration file: E:\nabbasi\data\nabbasi_web_Page\academic\my_courses\MATH127\final_project\Ww4_prob1.hmm

Enter observation sequence: 3462

Probability of the whole observation sequence: **0.000753852**

trellis diagram, Alpha (Forward)

Probabilities in logs			
S	:-2.4847066698	-4.0532425909	-5.6217785120
T	:-2.0794415417	-3.6479774628	-4.1588830834*

Probabilities

S	:+0.0833500000	+0.0173659725	+0.0036182004
T	:+0.1250000000	+0.0260437500	+0.0036182004

Column Sums of Probabilities

prob.	:+0.2083500000	+0.0434097225	+0.0434097225
-------	----------------	---------------	---------------

trellis diagram, Beta (backward)

Probabilities in logs			
S	:-5.6217785120	-4.0532425909	-1.56853
T	:-5.6217785120	-4.0532425909	-1.56853

Probabilities

S	:+0.0036182004	+0.0173659725	+0.20835
T	:+0.0036182004	+0.0173659725	+0.20835

Position : 1 2

trellis diagram, Gamma (Viterbi)

Probabilities in logs			
S	:-2.4847066698	-4.5641482115	-6.1562350000
T	:-2.0794415417*	-4.1588830834*	-4.1588830834*

Probabilities

S	:+0.0833500000	+0.0104187500	+0.0156250000
T	:+0.1250000000	+0.0156250000	+0.0156250000

Column Sums of Probabilities

prob.	:+0.2083500000	+0.0434097225	+0.0434097225
-------	----------------	---------------	---------------

Position : 1 2

State transition: T,T,S,T

Figure 14: SHMM MATLAB user interface

References

- [1] R.Durbin,S.Eddy,A.Krogh,G.Mitchison. Biological Sequence Analysis.
- [2] Dr Lior Pachter lecture notes. Math 127. Fall 2002. UC Berkeley.
- [3] Identification of genes in human genomic DNA. By Christopher Burge, Stanford University, 1997.