

POLITECNICO DI TORINO

Corso di Laurea Magistrale in Ingegneria Informatica

Tesi di Laurea Magistrale

Machine learning for big sequence data: Wavelet-compressed Hidden Markov Models



**POLITECNICO
DI TORINO**

Relatore

prof. Paolo Garza

Studente

Luca BELLO

Chalmers University of Technology

prof. Alexander Schliep

Giugno 2020

This work is subject to the Creative Commons Licence

LUCA BELLO

Department of Computer Science and Engineering

Chalmers University of Technology and University of Gothenburg

Abstract

Hidden Markov models are among the most important machine learning methods for the statistical analysis of sequential data, but they struggle when applied to big data. Their relative inefficiency has been addressed several times by the use of some compression techniques, either for the computation or for the data. This thesis explores the latter, with the application of a data compression technique based on wavelets and the subsequent adaptation of the main HMMs algorithms from the literature: the forward, Viterbi and Baum-Welch algorithms used to solve the evaluation, decoding and training problem respectively. The testing phase shows that this new technique generally yields equal or better results, obtaining some extremely high speedups in the training problem, making it even thousands of times faster; this enables the training of a HMM with big data on a commodity laptop.

Keywords: machine, learning, sequence, wavelet, compression, hidden, markov, models, viterbi, training.

Acknowledgements

Here I am, looking at this page and trying to write something that I can really call my own; something I can unmistakably recognize as marked with my personal style. Anyway, I will limit myself to just sprinkle some little magic here and there, sticking to a traditional acknowledgment format, because there are definitely people I need to thank. Also, I know for sure this will end up being too long, but I'm not the one reading, so the joke's on you.

The crazy series of events, life experiences and people that led me to this moment would be too complex to describe. I started university as a completely different person than I am now, moving out of my small hometown, away from my family, and joining a student dorm in a distant city. I wasn't prepared for what came after: the absurd number of sleepless nights devoted to study; questioning my life choices more times than I can remember; self-reflections that turned my soul inside out multiple times; abandoning everything to start over. And I definitely didn't know that this experience would have given me lifelong friendships with people that I love and respect. I'm incredibly grateful to them, and going through these five years without them would have been impossible. This is the moment I should stop talking about me and my exceptionally moving story and start talking about them. Before starting a foreseeable long list of people, I should thank my family, that has always been incredibly supportive throughout my academic career.

My dorm gave me a second family, and I need to thank every single person I met there since everyone contributed to enriching my personality with a small part of theirs, each in their own way. I want to thank Niccolò, for all the morning coffees, the nightly bitters, the foam swords duels and the flip flops thrown against the wall; Angelo, for his loyalty, the chicken curry, and for shouting my name out loud to all the girls he meets; Fabrizio, for being a constant inspiration and motivation to do better, for all the times he says "sorry", and for his kindness; and I want to thank the whole group "L'associazione" for all the gaming nights, the shared panic for the exams, the meals we prepared, the laughs, the tears, and all the time we spent and will spend together. I want to thank Luciana for all the tea we drank at night, for the penguin hugs, the talks sitting against the radiator, and for far more things than I could list. I want to thank Martina LL, my "spiritual sister", for being there when it mattered. I want to thank Emma for being each other's rambunctious sidekick during the swedish adventure, and for all the times she tried to drag me out of my room towards social activities.

There are tons of other people who would deserve some big thanks - and a cookie - for being special, but this page won't fit them all. To every single one of you: thanks, from the bottom of my heart.

Luca Bello, Gothenburg, June 2020

Contents

List of Figures	ix
1 Introduction	1
1.1 Motivation and previous research	1
1.2 Hidden Markov Models	3
1.3 Wavelet compression	5
2 Theory	9
2.1 Evaluation problem	9
2.1.1 Forward algorithm	10
2.1.2 Backward algorithm	11
2.2 Decoding problem	12
2.3 Training problem	13
2.3.1 Starting model	15
2.4 Compressed algorithms	15
2.5 Formal transformations	17
2.5.1 Forward algorithm	18
2.5.2 Backward algorithm	18
2.5.3 Viterbi algorithm	18
2.5.4 Baum-Welch algorithm	19
2.6 Relevant model parameters	20
2.6.1 State separation	20
2.6.2 Self-transition probabilities	22
3 Methods	23
3.1 Framework description	23
3.2 Implementation details	24
3.2.1 Parser	24
3.2.2 Numerical errors	24
3.2.3 State representation	25
3.2.4 Data generation	25
3.2.5 Wavelet compression	26
3.2.6 Logarithms summation	26
3.2.7 Saving results	27
3.2.8 Test automation	27
3.2.9 Other Python files	27

3.3	Standard algorithms	28
3.3.1	Evaluation problem	28
3.3.2	Decoding problem	28
3.3.3	Training problem	28
3.4	Compressed algorithms	29
3.4.1	Evaluation problem	30
3.4.2	Decoding problem	30
3.4.3	Training problem	30
4	Results	31
4.1	Testing setup	31
4.1.1	Choosing the parameters	31
4.1.2	Results evaluation	32
4.2	Evaluation problem	33
4.3	Decoding problem	34
4.4	Training problem	35
4.5	Circular topology	38
4.5.1	Evaluation problem	38
4.5.2	Decoding problem	39
4.5.3	Training problem	40
4.6	Left-to-right topology	42
4.6.1	Evaluation problem	42
4.6.2	Decoding problem	43
4.6.3	Training problem	44
5	Conclusion	47
5.1	Main takeaways	47
5.2	Wrapping up	47
	Bibliography	49

List of Figures

1.1	Collocation of this thesis in the context of different HMM types. . . .	2
1.2	A simple example of a Markov model: the state chain through time (top) and the graph representation of the model (bottom).	4
1.3	An example of the output of a hidden Markov model through some time steps. The state path Q is not observable, and O are the observations emitted in each state.	4
1.4	A signal decomposed with wavelets in its approximation and different levels of detail. Image taken from [1].	6
1.5	Typical sequence data undergoing block compression; the vertical red lines show the block borders in the case of an ideal compression. . . .	6
1.6	Example of block generation using a breakpoint array. During the query (thick red), when values are above the threshold (horizontal blue line), a breakpoint is returned (vertical blue line). This figure has been taken from [2].	7
2.1	Graphical representation of the lattice structure used for dynamic programming in the standard algorithms.	11
2.2	The structure on which the compressed algorithms operate; it generates from a subset of the links in the lattice structure.	16
2.3	Data generated using a two-state HMM with distributions $N(0, 1^2)$ and $N(10, 1^2)$	21
3.1	Structure of the implemented code (thin border) and its connection points to the external libraries used (thick border).	24
4.1	Relative difference between the $P(O \lambda)$ log probabilities of the compressed and standard algorithms. Since the errors are on a big negative log probability, the actual error magnitude is approximately of 10^{-5000}	33
4.2	Speedup on the evaluation problem using the compressed algorithm, including the input data processing time.	33
4.3	Relative difference between the fractions of errors in the estimated generating path.	34
4.4	Speedup on the decoding problem using the compressed algorithm, including the input data processing time.	34
4.5	Difference between the average KL-divergence for the compressed and standard algorithms.	35

4.6	Relative difference between the average error on the log transition probabilities of the compressed and standard algorithms.	36
4.7	Relative difference between the average error on the log probability of the starting state for the compressed and standard algorithms. . .	36
4.8	Speedup on the training problem using the compressed algorithm, including the input data processing time.	37
4.9	Relative difference between the $P(O \lambda)$ log probabilities of the compressed and standard algorithms.	38
4.10	Speedup on the evaluation problem using the compressed algorithm, including the input data processing time.	38
4.11	Relative difference between the fractions of errors in the estimated generating path.	39
4.12	Speedup on the decoding problem using the compressed algorithm, including the input data processing time.	39
4.13	Difference between the average KL-divergence for the compressed and standard algorithms.	40
4.14	Relative difference between the average error on the log transition probabilities of the compressed and standard algorithms.	40
4.15	Relative difference between the average error on the log probability of the starting state for the compressed and standard algorithms. . .	41
4.16	Speedup on the training problem using the compressed algorithm, including the input data processing time.	41
4.17	Relative difference between the $P(O \lambda)$ log probabilities of the compressed and standard algorithms.	42
4.18	Speedup on the evaluation problem using the compressed algorithm, including the input data processing time.	42
4.19	Relative difference between the fractions of errors in the estimated generating path.	43
4.20	Speedup on the decoding problem using the compressed algorithm, including the input data processing time.	43
4.21	Difference between the average KL-divergence for the compressed and standard algorithms.	44
4.22	Relative difference between the average error on the log transition probabilities of the compressed and standard algorithms.	44
4.23	Relative difference between the average error on the log probability of the starting state for the compressed and standard algorithms. . .	45
4.24	Speedup on the training problem using the compressed algorithm, including the input data processing time.	45

1

Introduction

Nowadays, many real-world problems are tackled trying to get insights from data with the application of machine learning techniques; it is easy to see that as the amount of data becomes massive, the algorithms necessarily need to be adapted and optimized towards a lower computational complexity to retain feasibility in reasonable execution time.

Hidden Markov Models (HMMs) are among the most important machine learning methods for the statistical analysis of sequential data. Their mathematical structure provides a wide range of applications and can serve as a basis for more complex models, making their study extremely meaningful. Some central and well-known algorithms for HMMs tend to struggle when applied on big data; this thesis aims to adapt them so they can be applied on a compressed representation of the data, preserving their usefulness and applicability.

1.1 Motivation and previous research

There are many examples of real-world applications showing the importance of hidden Markov models: computational finance [3], speech recognition [4], and more. Unfortunately, limitations have been encountered when working with big data due to the computational complexity of the standard algorithms. This problem has been tackled with many different approaches: one that has often proved its effectiveness is the idea of compression, to significantly reduce either the computations or the scale of the data.

To show some context and establish the relevance of this topic, some efforts and achievements will be described. For discrete-valued sequences, [5] shows how some speed improvements are obtained by text compression techniques, based on identifying repeated substrings in the observed input sequence and obtaining highly parallelizable algorithms. Bayesian computations that were often avoided in practice due to long running times have been accelerated in [6], showing considerable improvements. For continuous observations, Bayesian inference was tackled by [7] and [2] using wavelet compression; as summarized in Figure 1.1, what has been left uncovered are the possible speed improvements obtainable by applying compression techniques to the standard algorithms.

This thesis aims at covering that unexplored case, studying the effects of wavelet compression for regular HMMs and their possible impact on real-world problems. Specifically, I will make use of the concepts in Wiedenhoeft's papers [7] and [2], using part of the HaMMLET tool for wavelet compression; this is described more

	Standard	Bayesian
Discrete	Z.-U. M. Mozes S., Weimann O.	M. P. Mahmud, A. Schliep
Continuous	this thesis	S. A. Wiedenhoeft J., Brugel E.

Figure 1.1: Collocation of this thesis in the context of different HMM types.

in detail in Section 1.3.

1.2 Hidden Markov Models

This section provides basic information on hidden Markov models and introduces some needed notation, borrowed from the one used by Rabiner in [8]. Consider a system which can be described by a set of N distinct states, S_1, S_2, \dots, S_N , as shown in Figure 1.2. At regular discrete time intervals, the system goes through a state change and moves to another state (or to the same one) according to certain probabilities associated with the state; if the future evolution of the system depends solely on the current state, regardless of the system's history, this is a first-order Markov process.

Let the time be denoted with $t = 1, 2, \dots$, and the state of the system at a certain time with q_t . For a Markov process, it is true that:

$$P(q_t = S_j | q_{t-1} = S_i, q_{t-2} = S_k, \dots) = P(q_t = S_j | q_{t-1} = S_i). \quad (1.1)$$

The processes that are more often considered are *time-homogeneous*, meaning that the transition probabilities a_{ij} from a state to another are independent of time; these probabilities can be collected in a matrix A that obeys the standard stochastic constraints:

$$a_{ij} = P(q_t = S_j | q_{t-1} = S_i) \geq 0, \quad \text{and} \quad (1.2a)$$

$$\sum_{j=1}^N a_{ij} = 1 \quad (1.2b)$$

In the Markov model defined above and shown in Figure 1.2, each state corresponds to an observable event. That is not always the case; sometimes the observations are not the states of a Markov model, but they are related to the actual state chain which is *hidden*, or *not observable*. A hidden Markov model (as shown in Figure 1.3) has an underlying Markov model which works as described above, with the observations being a probabilistic function of the current state; this means that for any observation O_t

$$P(O_t | O_{t-1}, O_{t-2}, \dots, q_t, q_{t-1}, q_{t-2}, \dots) = P(O_t | q_t). \quad (1.3)$$

Based on the type of observations, Markov models can be divided in discrete, where the alphabet of observations is finite, and continuous, where it's not. Looking at the definition of a model, many interesting questions may arise: how likely is it that a sequence of observations has been generated by a certain model? What is the most likely state path corresponding to a sequence of observations? What is a good estimation of a model that suits the observations well? How can we apply Bayesian inference to train a model? The first three questions are, as defined by Rabiner in [8], the standard problems of HMMs named evaluation, decoding and training problem respectively. The aim of this thesis is to address these problems in the context of big data with a more efficient approach.

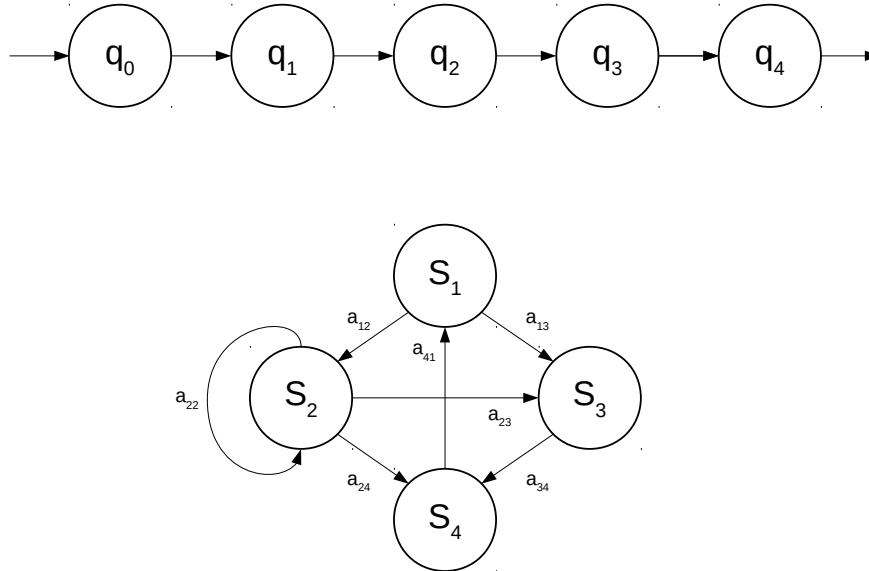


Figure 1.2: A simple example of a Markov model: the state chain through time (top) and the graph representation of the model (bottom).

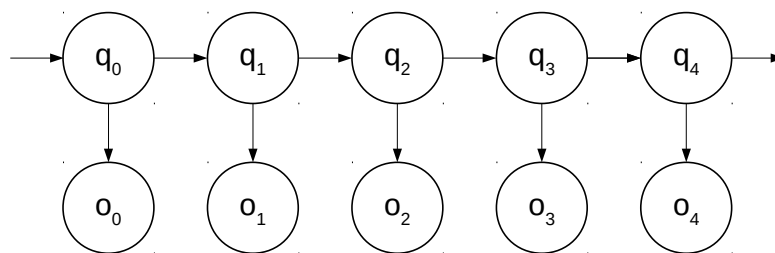


Figure 1.3: An example of the output of a hidden Markov model through some time steps. The state path Q is not observable, and O are the observations emitted in each state.

1.3 Wavelet compression

The compression technique that will be used throughout this thesis is based on wavelets. A wavelet is a function that resembles a small oscillation; the convolution of scaled versions of a wavelet with a signal of interest yields its wavelet transform. This is often used as a mathematical tool to analyze signals, obtain a different representation of them, and extract useful information. It is similar to the Fourier transform, with the main difference that the latter loses all information about the localization of a given frequency component.

A signal is usually decomposed using a certain wavelet and then, after some processing, often recomposed using the corresponding reconstruction wavelet. The decomposition yields an *approximation* of the signal and some *detail* coefficients, representing a smoothed version of the signal and the higher frequency information respectively. Many different wavelets can be used in a transform, based on the properties needed for a certain application:

- *size of support*, the interval where the wavelet is non-zero;
- *symmetry*, that influences the quality of localization;
- *number of vanishing moments*, a blindness to polynomials of a certain degree;
- *regularity*, also affecting frequency localization;
- *(bi-)orthogonality*, meaning that the decomposition and reconstruction wavelets form two distinct bases which are mutually orthogonal.

Figure 1.4 shows an example of signal decomposition using wavelets; this representation allows for effective de-noising against additive white Gaussian noise by simply setting to zero some detail coefficients above a certain threshold before reconstructing the signal. An interesting description of a possible use of wavelets is [9].

In the context of this thesis and following Wiedenhoeft's work in [7] and [2], the Haar wavelet is used to detect when a sequence of observations has a significant discontinuity in values, possibly indicating, under certain conditions, a change of state. An example result of the compression process is shown in Figure 1.5. Specifically, two main data structures from the HaMMLET tool will be used: the *breakpoint array* and the *integral array*. The former uses a certain threshold (discussed later with the implementation) to define a block structure by storing indexes of block delimiters, dividing the sequence of observations into blocks where the state can be considered to stay the same; the mechanism is shown in Figure 1.6. The latter contains sufficient statistics for each block, such as the sum of all elements inside it. Using two structures instead of one yields a more efficient implementation, as described in [2].

It is reasonable to expect that this compression method will work well when the states of the analyzed HMM are well-separated so that the block division is accurate, and when the self-loop probability for each state is high enough so that a single block contains more observations. In principle, this approach could also be applied to multivariate data, but that falls out of the scope of this thesis.

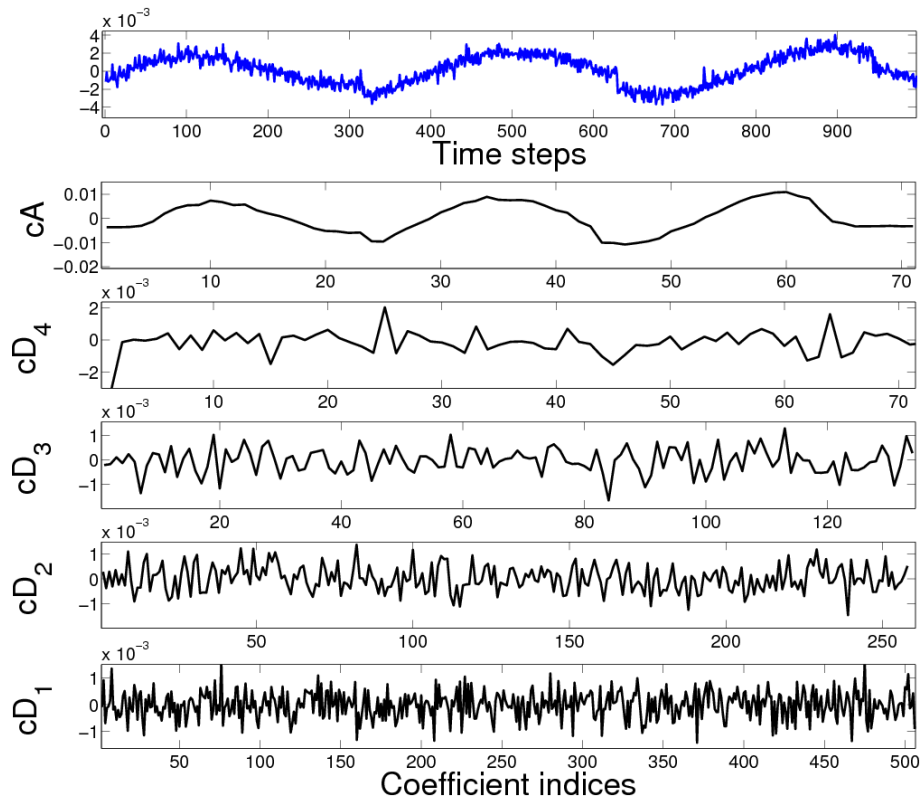


Figure 1.4: A signal decomposed with wavelets in its approximation and different levels of detail. Image taken from [1].

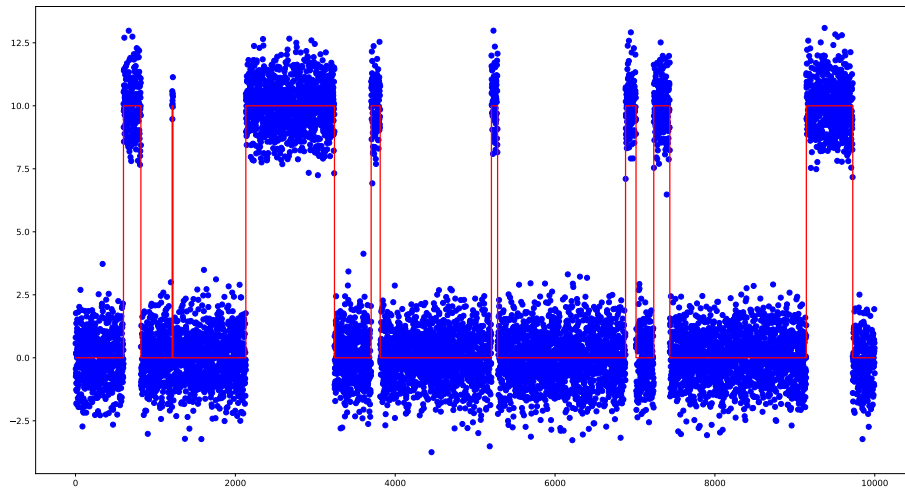


Figure 1.5: Typical sequence data undergoing block compression; the vertical red lines show the block borders in the case of an ideal compression.

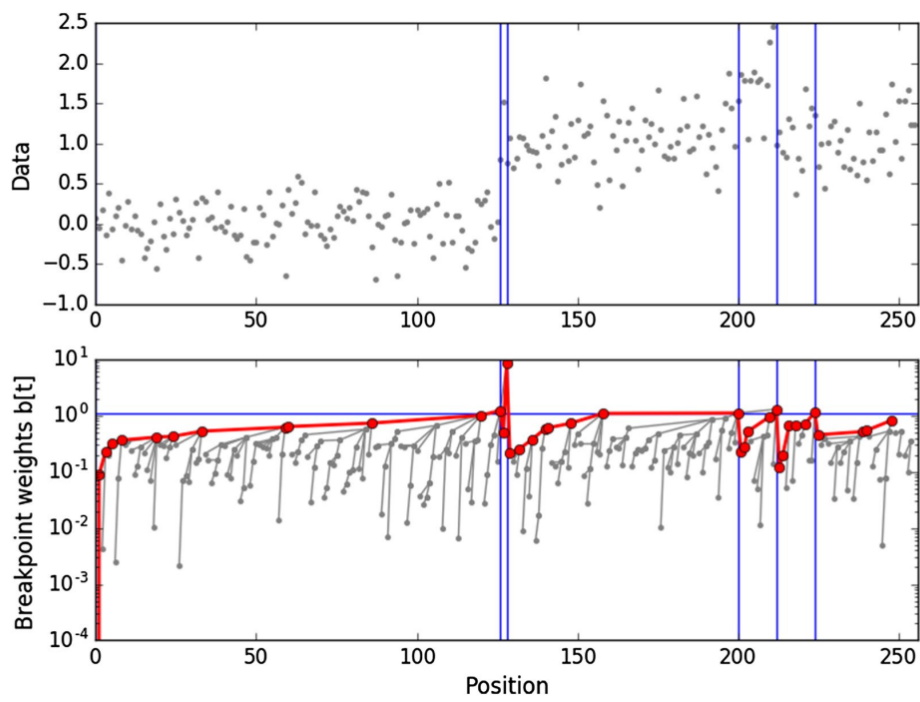


Figure 1.6: Example of block generation using a breakpoint array. During the query (thick red), when values are above the threshold (horizontal blue line), a breakpoint is returned (vertical blue line). This figure has been taken from [2].

2

Theory

This chapter tackles the theory behind hidden Markov models and their main problems of interest, using that as a stepping stone to analyze the theoretical transformations required by the data compression process. At this point the elements of a HMM should be clearly defined and denoted; again, the notation used is taken from Rabiner [8]:

- N , the number of states in the model; specifically, individual states will be denoted as $S = \{S_1, S_2, \dots, S_N\}$ and the state at time t will be denoted as q_t ;
- $A = \{a_{ij}\}$, the state transition probability distribution, where a_{ij} indicates the probability of going from state i to state j in one time step;
- $B = \{b_j\}$, the emission probability distribution of observations for each state j ;
- $\pi = \{\pi_i\}$, the initial state distribution.

To indicate the complete set of parameters, this compact notation is generally used,

$$\lambda = (N, A, B, \pi). \quad (2.1)$$

To add on this and clarify the notation used: the observation sequence is denoted by O , with the element observed at time t being indicated by O_t ; in the same way, the state sequence (also *state path*, *generating path*) is denoted by Q , and the state at time t is indicated by q_t . The HMMs that will be used in this thesis are continuous-valued; each state is associated and characterized by an emission distribution. For many applications, using Gaussian distributions is a good choice (e.g. autoregressive HMMs for speech recognition). As stated in the introduction, three main problems will be tackled; each one will be described and discussed in a separate section below.

2.1 Evaluation problem

The problem is about evaluating how well a specific sequence of observations is represented by a given model, through the computation of the probability that the observed sequence was produced by the model. Solving this problem is important because it allows the comparison of different models to decide which one better represents a sequence of observations. Formally, given the observation sequence $O = O_1, O_2, \dots, O_T$ and a model $\lambda = (A, B, \pi)$, the goal is to efficiently compute its likelihood. A very intuitive but inefficient way of doing it would be applying the law of total probability, enumerating every possible sequence of states and summing the

conditional probability of all observations over all those sequences:

$$P(O|\lambda) = \sum_{all\ Q} P(O|Q, \lambda) = \sum_{all\ Q} \pi_{q_1} b_{q_1}(O_1) a_{q_1 q_2} b_{q_2}(O_2) \dots a_{q_{T-1}} b_{q_T}(O_T) \quad (2.2)$$

For a specific state sequence, this equation starts in the state q_1 with probability π_{q_1} , produces the symbol O_1 with probability $b_{q_1}(O_1)$ and moves to the next state q_1 to q_2 with probability $a_{q_1 q_2}$, and then follows the same logic to cover the whole observation sequence. When looking at the computational complexity, the number of state paths obtained by enumeration is given by the dispositions with repetition of the states, that is N^T ; for each state, the number of calculations scales linearly with the length of the observation sequence; this yields a total computational complexity that is $\mathcal{O}(T \cdot N^T)$.

2.1.1 Forward algorithm

The standard algorithm used to solve the evaluation problem is the *forward* algorithm, which is much more efficient than the approach considered before. The key element is the forward variable $\alpha_t(i)$, defined as the joint probability of observing the sequence up to time t and being in state S_i at time t

$$\alpha_t(i) = P(O_1 O_2 \dots O_t, q_t = S_i | \lambda). \quad (2.3)$$

Through induction, the following procedure can be defined:

$$\alpha_1(i) = \pi_i b_i(O_1), \quad 1 \leq i \leq N \quad (2.4a)$$

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(O_{t+1}), \quad 1 \leq t \leq T-1, \quad 1 \leq j \leq N \quad (2.4b)$$

$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i). \quad (2.4c)$$

The first step is the initialization of the forward variable in (2.4a) using the initial probability distribution π . The actual induction happens on the second step (2.4b), where α_{t+1} is calculated using the forward variables at the previous instant α_t ; to make that calculation for a certain state S_j , it is necessary to consider the probability of getting there from any state S_i multiplied by the probability of observing the symbol O_{t+1} . This procedure terminates with (2.4c), that simply takes the sum of the forward variables over all the states. The procedure can be visualized well through Figure 2.1.

Looking at the computational complexity, the number of calculations for each observation is N^2 (N per each state); repeating this for the whole sequence length gives a complexity that is $\mathcal{O}(T \cdot N^2)$, which is faster than the previous approach, especially for increasingly long sequences.

Very closely tied to this algorithm is the backward algorithm. Although it's not necessary to solve the evaluation problem, it can be helpful in the solution of both the decoding and the training problems.

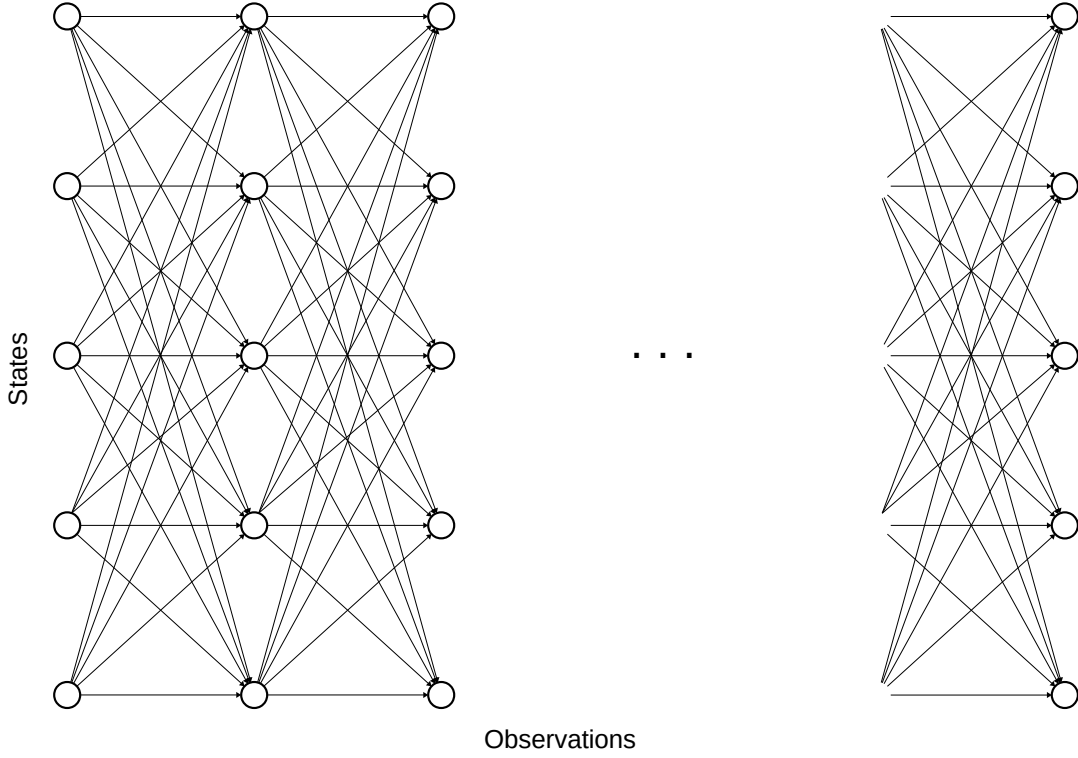


Figure 2.1: Graphical representation of the lattice structure used for dynamic programming in the standard algorithms.

2.1.2 Backward algorithm

The backward variable is defined as the probability of the partial observation sequence from $t + 1$ to the end, given a certain state at time t and the model

$$\beta_t(i) = P(O_{t+1}O_{t+2} \dots O_T | q_t = S_i, \lambda). \quad (2.5)$$

This variable can also be calculated using the lattice structure shown in Figure 2.1, in an inductive fashion with the use of dynamic programming:

$$\beta_T(i) = 1, \quad 1 \leq i \leq N, \quad \text{and} \quad (2.6a)$$

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \beta_{t+1}(j), \quad t = T-1, T-2, \dots, 1, \quad 1 \leq i \leq N. \quad (2.6b)$$

The initialization step in (2.6a) chooses an arbitrary starting point; the induction (2.6b) computes the other variables by accounting for all the possible states the system could have been in at the following time step. The calculation of the computational complexity of this algorithm follows the same reasoning of the forward algorithm; this leads to conclude that the backward algorithm also has a complexity that is $\mathcal{O}(T \cdot N^2)$.

2.2 Decoding problem

A very common situation when dealing with hidden Markov models is the need to predict the generating state path of a certain observed sequence. There is no exact answer to this question; the goal is to find the solution that better fits the data. Many optimization criteria can be used based on the definition of *better*. An example is finding the most likely state individually for each observation; using this criterion, the key variable is

$$\gamma_t(i) = P(q_t = S_i | O, \lambda), \quad (2.7)$$

which is the probability of being in state S_i at time t , given both the observation sequence and the model. This variable can be expressed using the forward and backward variables in the following way

$$\gamma_t(i) = \frac{\alpha_t(i)\beta_t(i)}{P(O|\lambda)} = \frac{\alpha_t(i)\beta_t(i)}{\sum_{j=1}^N \alpha_t(j)\beta_t(j)}, \quad (2.8)$$

because the forward variable accounts for the observations before time t , while the backward variable refers to the observations after t . The most likely state at time t is easily obtained through the γ_t variable by looking at which state has the highest associated probability:

$$q_t = \operatorname{argmax}_{1 \leq i \leq N} [\gamma_t(i)], \quad 1 \leq t \leq T. \quad (2.9)$$

Although this approach obtains the highest number of expected correct matches between predicted and actual state, it is often discarded since it disregards the transition probabilities of the model. In particular, if some transitions have zero probability, meaning they cannot happen, this criterion would still be able to include them in the result.

For this reason, it is often more interesting to compute the most likely generating path; mathematically, this means finding the maximization of $P(Q|O, \lambda)$. The standard algorithm used to solve this problem is the *Viterbi* algorithm; as the forward algorithm, it uses a dynamic programming approach that allows to reduce computational complexity. The Viterbi algorithm substitutes the summations of the forward algorithm with a maximization; the key variable to calculate is $\delta_t(i)$, defined as the highest probability along a single path, after the first t observations and ending in state i :

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_1 q_2 \dots q_t = i, O_1 O_2 \dots O_t | \lambda) \quad (2.10)$$

The most likely path will be the argument of this maximization over all the states considering the whole observations sequence; it can be defined in the notation as ψ . Through induction it is possible to write the following equations:

$$\begin{aligned} \delta_1(i) &= \pi_i b_i(O_1), & 1 \leq i \leq N \\ \psi_1 &= 0 \end{aligned} \quad (2.11a)$$

$$\delta_t(j) = \max_{1 \leq i \leq N} \left[\delta_{t-1}(i) a_{ij} \right] b_j(O_t), \quad 2 \leq t \leq T, \quad 1 \leq j \leq N \quad (2.11b)$$

$$\psi_t(j) = \operatorname{argmax}_{1 \leq i \leq N} \left[\delta_{t-1}(i) a_{ij} \right], \quad 2 \leq t \leq T, \quad 1 \leq j \leq N$$

$$P^* = \max_{1 \leq i \leq N} \left[\delta_T(i) \right] \quad (2.11c)$$

$$q_T^* = \operatorname{argmax}_{1 \leq i \leq N} \left[\delta_T(i) \right]$$

$$q_t^* = \psi_{t+1}(q_{t+1}^*), \quad t = T-1, T-2, \dots, 1 \quad (2.11d)$$

The initialization phase in (2.11a) starts with an empty solution for the state ψ_1 . The induction in (2.11b) relies on the lattice structured defined previously in Figure 2.1; the contribution of the previous δ_t variables is given through the max operator instead of using summation. During this process, the argmax is saved to later use the variable ψ to recover the most likely path. The termination (2.11c) happens at the end of the observations sequence. After the final step, the most likely state path associated with the sequence is given by backtracking through the ψ variables as shown in (2.11d).

To discuss the computational complexity of this algorithm, the same points of the forward algorithm can be made, leading to affirm that the complexity of the Viterbi algorithm is also $\mathcal{O}(T \cdot N^2)$.

2.3 Training problem

Real-world applications present many scenarios where the HMMs' parameters are not explicitly known. The relevance of this problem is easily shown by noting that a model gives a lot of insights on the system; moreover, its accuracy was a key assumption in the previous computations. The goal is to find the model λ that maximizes the probability of the observation sequence given the model, $P(O|\lambda)$. Unfortunately, this problem is very complex and there is no known way to analytically solve this maximization problem for any given finite observation sequence.

Despite this, a number of techniques can be used to locally maximize $P(O|\lambda)$; a very popular one is the *Baum-Welch* method, that starts from a guess of the model and iteratively performs reestimations of the parameters to improve it. This algorithm introduces a new key variable: $\xi(i, j)$, the probability of being in state S_i at time t and in state S_j at time $t+1$

$$\xi_t(i, j) = P(q_t = S_i, q_{t+1} = S_j | O, \lambda). \quad (2.12)$$

It can be useful to express this equation using the forward and backward variables. In fact, the forward variable $\alpha_t(i)$ accounts for the observations from the first one up to O_t in state S_i ; the backward variable $\beta_{t+1}(j)$ does the complementary job, considering the observation sequence starting in state S_j and from observation O_{t+1} up to the last one. The step between t and $t+1$ has been left out: to tie the two

variables, it is necessary to include the probability of transitioning from state S_i to S_j and observing O_{t+1} , which is $a_{ij}b_j(O_{t+1})$. The new formulation of ξ_t can be written as

$$\xi_t(i, j) = \frac{\alpha_t(i)a_{ij}b_j(O_{t+1})\beta_{t+1}(j)}{P(O|\lambda)} = \frac{\alpha_t(i)a_{ij}b_j(O_{t+1})\beta_{t+1}(j)}{\sum_{p=1}^N \sum_{q=1}^N \alpha_t(p)a_{pq}b_q(O_{t+1})\beta_{t+1}(q)}. \quad (2.13)$$

By looking at the definition of the γ_t variable given in (2.7), it can be related with ξ_t :

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j) \quad (2.14)$$

Recalling the previous definition of γ_t given in (2.7) is important to notice that by summing $\gamma_t(i)$ over t , the obtained quantity can be interpreted as the expected number of times that the state S_i is visited, or equivalently as the expected number of transitions from state S_i (if we exclude the last observation at time T):

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{expected number of transitions from state } S_i \quad (2.15)$$

In a similar way, the sum of $\xi_t(i, j)$ over t can be interpreted as the expected number of transitions from S_i to S_j :

$$\sum_{t=1}^{T-1} \xi_t(i, j) = \text{expected number of transitions from state } S_i \text{ to } S_j \quad (2.16)$$

These interpretations lead to the definition of two reestimation formulas for the initial distribution and the transition probabilities:

$$\bar{\pi}_i = \gamma_1(i) \quad (2.17a)$$

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} \quad (2.17b)$$

The HMMs that have been considered throughout this thesis work have continuous emission densities, that in the most general case can be written as:

$$b_j(O) = \sum_{m=1}^M c_{jm} \mathfrak{R}[O, \mu_{jm}, U_{jm}] \quad (2.18)$$

where O is the observations sequence, c_{jm} is the mixture coefficient of the m -th mixture in state S_j and \mathfrak{R} is a log-concave or elliptically symmetric density with mean vector μ_{jm} and covariance matrix U_{jm} , again for the m -th mixture in state S_j . It can be shown ([10–12]) that the reestimation formulas for the coefficients of the mixture density have the following form

$$\bar{c}_{jk} = \frac{\sum_{t=1}^T \gamma_t(j, k)}{\sum_{t=1}^T \sum_{k=1}^M \gamma_t(j, k)}, \quad (2.19a)$$

$$\bar{\mu}_{jk} = \frac{\sum_{t=1}^T \gamma_t(j, k) \cdot O_t}{\sum_{t=1}^T \gamma_t(j, k)}, \quad \text{and} \quad (2.19b)$$

$$\bar{U}_{jk} = \frac{\sum_{t=1}^T \gamma_t(j, k) \cdot (O_t - \mu_{jk})(O_t - \mu_{jk})'}{\sum_{t=1}^T \gamma_t(j, k)}. \quad (2.19c)$$

where prime denotes the transposition of the vector and $\gamma_t(j, k)$ is simply $\gamma_t(j)$ relative to the k -th mixture component. According to the context and the scope of the thesis, the mixture model is reduced to a univariate Gaussian distribution; the emission density for a state S_j can be rewritten as

$$b_j(O) = \mathfrak{R}[O, \mu_j, \sigma^2], \quad (2.20)$$

where μ_j is the mean and σ^2 is the variance of the Gaussian distribution associated with the state S_j . Thus, the reestimation formulas can also be simplified (e.g. by getting rid of the mixture weight coefficients) and rewritten:

$$\bar{\mu}_j = \frac{\sum_{t=1}^T \gamma_t(j) \cdot O_t}{\sum_{t=1}^T \gamma_t(j)} \quad (2.21a)$$

$$\bar{\sigma}^2_j = \frac{\sum_{t=1}^T \gamma_t(j) \cdot (O_t - \mu_j)^2}{\sum_{t=1}^T \gamma_t(j)} \quad (2.21b)$$

Applying the reestimation formulas (2.17a), (2.17b), (2.21a) and (2.21b) produces a reestimated model $\bar{\lambda}$; the Baum-Welch algorithm guarantees that either the original model λ is a critical point of the likelihood function (the result would be $\lambda = \bar{\lambda}$) or the model $\bar{\lambda}$ is more likely than the previous one, meaning that $P(O|\bar{\lambda}) > P(O|\lambda)$. The iteration of this procedure converges to a local maximum and produces a maximum likelihood estimate of the model, providing a solution to the training problem.

2.3.1 Starting model

The Baum-Welch procedure requires the definition of a starting model; even though the number of states is generally known (or can be guessed or estimated), to obtain good results a good definition of the starting model $\lambda = (A, B, \pi)$ is necessary. Unfortunately, most of the time little knowledge is possessed about the system; thus, there is no straightforward answer to this problem. As discussed by Rabiner in [8], experience shows that for A and π either random or uniform initial estimates are adequate for useful parameters' reestimation. For continuous emission distributions B , the starting parameters are essential. Such parameters can be obtained with several techniques, such as manual segmentation of the observation sequence into states or k-means segmentation with clustering.

2.4 Compressed algorithms

The wavelet compression of the data has a big impact on the mechanisms of the algorithms. The compression removes the necessity to consider the observations

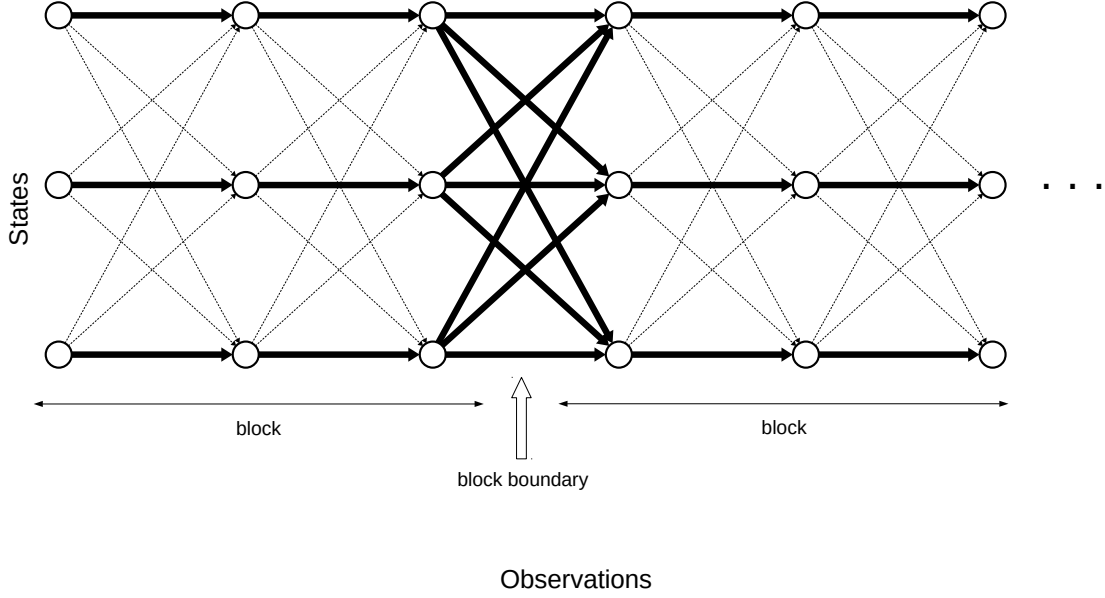


Figure 2.2: The structure on which the compressed algorithms operate; it generates from a subset of the links in the lattice structure.

individually, allowing to focus on groups of them called *blocks*. As anticipated in Section 1.3, a block is a series of observations where the underlying state can be considered constant and with sufficient statistics to perform the computations required in the algorithms of interest. The structure used by all the algorithms above, as said many times, is the one in Figure 2.1; not knowing anything about the generating path of the observations forces to consider all the possible ones, evaluating at each time step the contribute of every possible state. Compressing the observations opens up new possibilities by modifying that lattice structure into the one represented in Figure 2.2. Since inside a block the state can be assumed to remain the same, the transitions from other states are reputed too unlikely and thus ignored.

To talk formally about compression and its impact on the computations, it is necessary to introduce some notation; this is done taking [13] as a starting point, but applying some changes to avoid conflicts with already defined symbols and to put more emphasis on some concepts.

A partition of the observations in blocks can be denoted as $Y := \{Y_w\}_{w=1}^W$, where Y_w is a single block and W is the number of blocks forming the partition. A block contains n_w observations; each one is referred to using the symbol $y_{w,k}$ which is indexed by the block number w and by the position inside the block k . The summary statistics gathered for each block are the following;

$$n_w, \quad \Sigma_{1,w} := \sum_{k=1}^{n_w} y_{w,k}, \quad \Sigma_{2,w} := \sum_{k=1}^{n_w} y_{w,k}^2 \quad (2.22)$$

To figure out how the computation varies with the introduction of this block structure, it is useful to thoroughly examine the calculations in Wiedenhoeft's PhD thesis [13] where HaMMLET was developed. Remembering the definition of the forward variable at (2.3) and its computation formula at (2.4b), being inside a block Y_w only allows self-transitions; this reduces the computation of the next forward variable to

$$\alpha_{t+1}(j) = \alpha_t(j) a_{jj} b_j(O_{t+1}). \quad (2.23)$$

Using induction, the forward variable relative to the whole block has to account for $n_w - 1$ self-transitions, one transition to state S_j and n_w emissions; it can be expressed as

$$\alpha_w(j) = \sum_{i=1}^N \left[\alpha_{w-1}(i) a_{ij} \right] a_{jj}^{n_w-1} \prod_{k=1}^{n_w} b_j(y_{w,k}). \quad (2.24)$$

An analogous point can be made on the backward and the Viterbi algorithms; the key part, though, is that this formula still relies on individual observations. To exploit summary blocks statistics, the term accounting emissions and self-transitions within a block can be rewritten by making the Gaussian emissions explicit

$$a_{jj}^{n_w-1} \prod_{k=1}^{n_w} b_j(y_{w,k}) = \frac{a_{jj}^{n_w-1}}{\sqrt{2\pi}^{n_w} \sigma_j^{n_w}} \exp \left(- \sum_{k=1}^{n_w-1} \frac{(y_{w,k} - \mu_j)^2}{2\sigma_j^2} \right). \quad (2.25)$$

The factors outside the exponential can be brought in, also providing implementation advantages discussed later in Section 3.2.2. This yields

$$\exp \left(- \sum_{k=1}^{n_w-1} \frac{(y_{w,k} - \mu_j)^2}{2\sigma_j^2} + (n_w - 1) \log(a_{jj}) - n_w \log(\sigma_j) - n_w \log(\sqrt{2\pi}) \right). \quad (2.26)$$

The exponent can finally be rewritten using the blocks summary statistics:

$$E_w(j) := \frac{2\mu_j \Sigma_{1,w} - \Sigma_{2,w}}{2\sigma_j^2} + K(n_w, j), \quad \text{and} \quad (2.27a)$$

$$K(n_w, j) := (n_w - 1) \log(a_{jj}) - n_w \left(\log(\sigma_j) + \frac{\mu_j^2}{2\sigma_j^2} + \frac{1}{2} \log(2\pi) \right). \quad (2.27b)$$

As pointed out in Wiedenhoeft's PhD thesis [13], an equivalent term can be easily derived also for non-Gaussian emissions that belong to the exponential distribution family.

2.5 Formal transformations

To perform the other calculations, the equations have to be adapted using the reformulation above. The following sections contain the adaptation of the algorithms to the compression scheme.

2.5.1 Forward algorithm

Restructuring the forward algorithm does not bring a new meaning to the new variable; for a block, $\alpha_w(i)$ is the approximation of the uncompressed forward variable at the end of the block. The induction phase has already been defined in (2.24); adding the other steps yields:

$$\alpha_1(i) = \pi_i e^{E_1(i)}, \quad 1 \leq i \leq N \quad (2.28a)$$

$$\alpha_w(j) = \left[\sum_{i=1}^N \alpha_{w-1}(i) a_{ij} \right] e^{E_w(j)}, \quad 1 \leq w \leq W, \quad 1 \leq j \leq N \quad (2.28b)$$

$$P_Y(O|\lambda) = \sum_{i=1}^N \alpha_W(i). \quad (2.28c)$$

2.5.2 Backward algorithm

The backward algorithm follows a very similar transformation; $\beta_w(i)$ is defined here as the backward variable at the start of a block:

$$\beta_W(i) = 1, \quad 1 \leq i \leq N \quad (2.29a)$$

$$\beta_w(i) = \sum_{j=1}^N a_{ij} e^{E_{w+1}(j)} \beta_{w+1}(j), \quad w = W-1, W-2, \dots, 1, \quad 1 \leq i \leq N. \quad (2.29b)$$

2.5.3 Viterbi algorithm

The Viterbi algorithm is based on the forward algorithm with the substitution of the sum over all the states with the max operator; again, all the variables refer to the end of a block:

$$\begin{aligned} \delta_1(i) &= \pi_i e^{E_1(i)}, \quad 1 \leq i \leq N \\ \psi_1 &= 0, \end{aligned} \quad (2.30a)$$

$$\delta_w(j) = \max_{1 \leq i \leq N} \left[\delta_{w-1}(i) a_{ij} \right] e^{E_w(j)}, \quad 1 \leq w \leq W, \quad 1 \leq j \leq N \quad (2.30b)$$

$$\psi_t(j) = \operatorname{argmax}_{1 \leq i \leq N} \left[\delta_{w-1}(i) a_{ij} \right], \quad 1 \leq w \leq W, \quad 1 \leq j \leq N$$

$$P^* = \max_{1 \leq i \leq N} \left[\delta_W(i) \right], \quad (2.30c)$$

$$q_W^* = \operatorname{argmax}_{1 \leq i \leq N} \left[\delta_W(i) \right],$$

$$q_w^* = \psi_{w+1}(q_{w+1}^*), \quad w = W-1, W-2, \dots, 1. \quad (2.30d)$$

2.5.4 Baum-Welch algorithm

The Baum-Welch algorithm is more complex than the others, having more variables to calculate for the parameters reestimation. Following the same order of Section 2.3, the first computation of interest is $\xi_t(i, j)$, as defined in (2.13). To reason about the following computations, it is important to remember that both the compressed forward and backward variables refer to the end of a block. Given this, two different situations happen based on computing ξ inside or outside a block. Inside a block, it is easy to see that $\xi_t(i, j) = 0$ for $i \neq j$; if w indicates the block, when $i = j$ it becomes:

$$\begin{aligned} \xi_t(i, i) &= \frac{\alpha_t(i) a_{ii} b_i(O_{t+1}) \beta_{t+1}(i)}{P(O|\lambda)} = \\ &= \frac{\alpha_w(i)}{a_{ii}^{n_w-t} \prod_{k=t}^{n_w} b_i(O_k)} a_{ii} b_i(O_{t+1}) \beta_w(i) a_{ii}^{n_w-(t+1)} \prod_{k=t+2}^{n_w} b_i(O_k) \frac{1}{P_Y(O|\lambda)} = \\ &= \frac{\alpha_w(i) \beta_w(i)}{P_Y(O|\lambda)} \end{aligned} \quad (2.31)$$

Since the right expression is not dependent on t , inside a block the variable is constant over time. Instead, at the boundary between two blocks:

$$\begin{aligned} \xi_t(i, j) &= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{P(O|\lambda)} = \\ &= \alpha_w(i) a_{ij} b_j(O_{t+1}) \beta_{w+1}(j) a_{jj}^{n_{w+1}-1} \prod_{k=t+2}^{n_{w+1}} b_j(O_k) \frac{1}{P_Y(O|\lambda)} = \\ &= \frac{\alpha_w(i) a_{ij} e^{E_{w+1}(j)} \beta_{w+1}(j)}{P_Y(O|\lambda)} \end{aligned} \quad (2.32)$$

For the purpose of rewriting reestimation equations, it is useful to define the ξ variable for a block in the following way:

$$\begin{aligned} \xi_w(i, j) &= \sum_{t \in Y_w} \xi_t(i, j) = \\ &= \frac{1}{P(O|\lambda)} \cdot \begin{cases} (n_w - 1) \alpha_w(i) \beta_w(i) + \alpha_w(i) a_{ij} e^{E_{w+1}(j)} \beta_{w+1}(j) , & \text{for } i = j \wedge w \neq W \\ (n_W - 1) \alpha_W(i) \beta_W(i) , & \text{for } i = j \wedge w = W \\ \alpha_w(i) a_{ij} e^{E_{w+1}(j)} \beta_{w+1}(j) , & \text{for } i \neq j \wedge w \neq W \\ 0 , & \text{for } i \neq j \wedge w = W \end{cases} \end{aligned} \quad (2.33)$$

Moving forward, it is interesting to note that by interpreting $\gamma_t(i)$ as the probability of visiting the state S_i at time t , the variable is also constant over t inside a block (also implied from the result above); this means that for any t inside a block, any $\gamma_t(i)$ can be representative for the whole block; recalling that the forward and backward variables both refer to the end of a block, equation (2.8) can be rewritten:

$$\gamma_t(i) = \frac{\alpha_t(i) \beta_t(i)}{P(O|\lambda)} = \frac{\alpha_w(i) \beta_w(i)}{P_Y(O|\lambda)} \quad (2.35)$$

It is worth noting that this reformulation correctly maintains the definition given in (2.14). For convenience it is useful to define $\gamma_w(i)$ as the representative value for a block, which means that for all the t associated with a block, $\gamma_w(i) = \gamma_t(i)$.

The reestimation formulas at (2.17a), (2.17b), (2.21a) and (2.21b) follow the new definitions of $\xi_t(i, j)$ and $\gamma_t(i)$. Using the equations above, they can be rewritten:

$$\bar{\pi}_i = \gamma_1(i) \quad (2.36a)$$

$$\bar{a}_{ij} = \frac{\sum_{w=1}^W \xi_w(i, j)}{\sum_{w=1}^W [n_w \gamma_w(i)] - \gamma_W(i)} \quad (2.36b)$$

The mean and standard deviation reestimations follow a slightly more complex reformulation, both for the general mixture and the single Gaussian distribution. In particular, (2.21a) multiplies the single observation value by the respective $\gamma_t(i)$. Since $\gamma_t(i)$ is constant inside a block, the equation can be rewritten as:

$$\bar{\mu}_j = \frac{\sum_{t=1}^T \gamma_t(j) O_t}{\sum_{t=1}^T \gamma_t(j)} = \frac{\sum_{w=1}^W \gamma_w(j) \sum_{k=1}^{n_w} y_{w,k}}{\sum_{w=1}^W \gamma_w(j) \cdot n_w} = \frac{\sum_{w=1}^W \gamma_w(j) \cdot \Sigma_{1,w}}{\sum_{w=1}^W \gamma_w(j) \cdot n_w} \quad (2.37)$$

The same reasoning applies to the variance reestimation:

$$\bar{\sigma}_j^2 = \frac{\sum_{t=1}^T \gamma_t(j) (O_t - \mu_j)^2}{\sum_{t=1}^T \gamma_t(j)} = \frac{\sum_{w=1}^W \gamma_w(j) [\Sigma_{2,w} - 2\bar{\mu}_j \Sigma_{1,w} + n_w \bar{\mu}_j^2]}{\sum_{w=1}^W \gamma_w(j) \cdot n_w} \quad (2.38)$$

2.6 Relevant model parameters

The computational complexity of the compressed algorithms is lower than the standard versions; since the blocks are used instead of the individual observations, it goes from $\mathcal{O}(T \cdot N^2)$ to $\mathcal{O}(W \cdot N^2)$, where W is the total number of blocks obtained from the compression. The complexity analysis refers to an infinite amount of data; the actual efficiency gain depends on several factors: the separation of the states, the self-transition probabilities, and the implementation details (which will be discussed in Chapter 3).

These parameters also impact the results errors due to the approximation introduced by the block compression (as in (2.24)). Trying to find some conditions for which the compressed algorithms always work well (or badly) in terms of speed and accuracy is one of the main goals of this thesis, leading to a sensible choice on which set of algorithms to utilize in different situations.

This section presents the most relevant factors that influence the quality of the results, discussing their relevance and eventual conditions under which the compressed algorithms should yield a substantial advantage over the standard ones.

2.6.1 State separation

The state separation is a key factor in the compressed algorithms; if the states are well-separated, it is easier to distinguish between two of them. This means that the

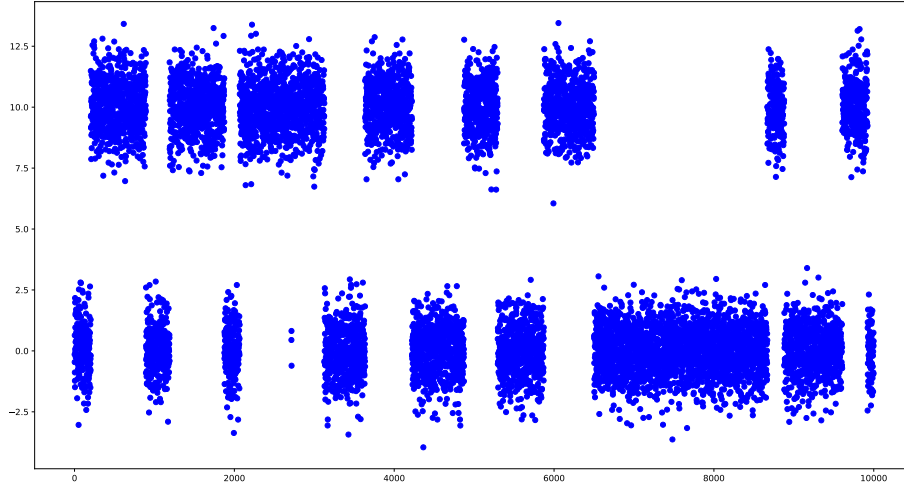


Figure 2.3: Data generated using a two-state HMM with distributions $N(0, 1^2)$ and $N(10, 1^2)$.

wavelet compression will produce blocks with clear boundaries, because some detail coefficients of the Maxlet transform will be high due to a bigger jump between two observations values belonging to different states. Figure 2.3 shows an example of data with clear distinction between the states.

If the states are well-separated, the approximation made by the block compression of neglecting some states' contribution becomes more accurate. This is true because the emission probability of an observation for a state that did not generate it gets very close to zero. In this context, it is useful to formally define a measure η_{S_1, S_2} of how well-separated two states are. Since for a Gaussian distribution $N(\mu, \sigma^2)$ it is known that 99.97% of the values lie within three standard deviations from the mean, a good measure definition could be:

$$\eta_{S_1, S_2} = \frac{|\mu_{S_1} - \mu_{S_2}|}{3(\sigma_{S_1} + \sigma_{S_2})} \quad (2.39)$$

The power of this definition lies in the fact that it is easy to understand visually: if $\eta_{S_1, S_2} = 1$ it means that the two distributions touch exactly after their respective 3σ ; if $\eta_{S_1, S_2} \ll 1$, the distributions overlap for a significant portion; if $\eta_{S_1, S_2} \gg 1$, then the distributions are clearly separated. Also, a measure of $\eta_{S_1, S_2} = 0$ indicates that the two distributions have the same mean. When having more than two states, this separation should apply between every pair of consecutive distributions (ordered by mean). This η will be denoted as *separation coefficient* for easier reference.

It should be clear that a higher state separation should reduce the error between the compressed and the standard algorithms; a more precise analysis will be conducted with the evaluation of the results in Chapter 4.

2.6.2 Self-transition probabilities

High self-transition probabilities should help the compression process in more than one way: for example, they make the blocks bigger, directly implying the production of fewer blocks for a fixed amount of observations and thus saving precious computation time. The impact on accuracy is hard to evaluate: for a fixed observation sequence, having fewer blocks could increase the error because more state paths are ignored; on the other hand, the contribution of the other state paths becomes smaller, thus decreasing the approximation error. To better understand the actual effects of this parameter, an extensive testing process will be discussed in Chapter 4.

3

Methods

This chapter describes the implementation structure, covering in detail the differences between the standard algorithms and their compressed version. Moreover, typical problems (e.g. numerical precision) will be addressed, discussing applied solutions and possible improvements.

The language of choice is C++, after considering others such as Python and R; it sacrifices some simplicity for the sake of efficiency, both in terms of speed and memory management. A great and popular compromise is to expose some Python bindings to an internal C++ structure, providing a simpler interface to the underlying complex implementation.

3.1 Framework description

In the design phase of a new project it is very useful to draw a scheme describing the class structure as the one in Figure 3.1, also listing what external tools will be used and how they will interface to the main components. To have more control over the code and allow for fair speed comparisons, both the HMM representation and the standard algorithms have been implemented from scratch.

The external components used are: CXXopts, to parse input arguments from command line; HaMMLET [13], for wavelet compression and the related data structures; Pomegranate, a Python framework used for data generation.

The tool has been named WaHMM, after **W**avelets **H**idden **M**arkov **M**odel and following the style of HaMMLET. The core interfaces with the external components through the **parser** and the **Compressor** elements. The standard algorithms and their compressed version have been separated in different files for easier management.

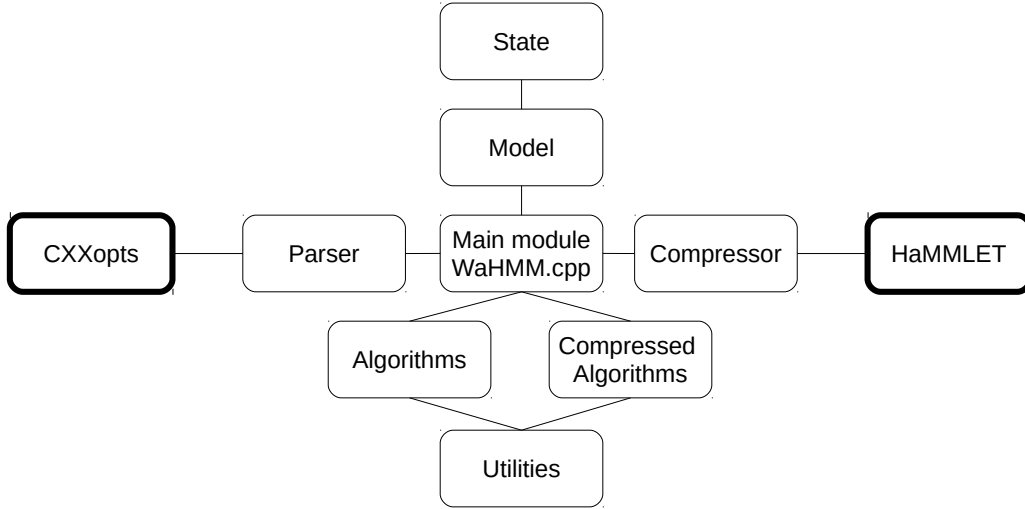


Figure 3.1: Structure of the implemented code (thin border) and its connection points to the external libraries used (thick border).

3.2 Implementation details

This section will describe several implementation choices, the reasoning behind them, and their impact on the results.

3.2.1 Parser

The parser allows for easy interaction with the tool through the specifications of various options using the common UNIX style. Although it can be quite verbose, it is a very effective instrument to define input parameters and which algorithms to execute. This interface makes it possible to input a model through command line or file, as well as giving input observations as a space-separated list of floats in a file or with a binary file format. Several options allow the user to choose which algorithms to execute, together with controls for verbosity and saving the results to files.

3.2.2 Numerical errors

Any kind of numerical method or scientific computation faces the problem of numerical errors. Representing real numbers on a machine is one thing, but observing that representation forces their decimal expansion to be truncated at some point. There is an entire sub-field of programming language theory called “exact real-number com-

putation” devoted to representing real numbers on computers; more information on that can be found on Haskell’s wiki page [14]. At the end of the day, the smallest difference between two numbers that a computer can recognize is called *machine epsilon*; if their difference is smaller, it is rounded to zero producing a rounding error. The machine epsilon is platform-dependent, but generally, it can be close to 10^{-16} . Since the product of probabilities can get very small fairly quickly, this problem would cause increasingly bigger errors on all computations (e.g. forward matrix after some time steps).

To address this well-known problem, many approaches are possible: to achieve higher precision, a double data type is used to store results, redefined as `wahmm::real_t`; the workspace `wahmm` is used to avoid naming conflicts with libraries such as HaMM-LET. For numerical problems, the strategy followed in this thesis is to use logarithms of probabilities. There are several advantages with this approach: the logarithm naturally scales the $[0, 1]$ interval to $(-\infty, 0]$ so that when probabilities get smaller, their logarithm becomes more negative; all the products required by the algorithms become summations, which is much easier and faster to perform. Using a logarithmic space also saves computations in the compressed algorithms; in fact, this transformation removes the exponential function in (2.27a), making it a natural choice for the use in these algorithms.

3.2.3 State representation

For the scope of this thesis, a state is associated with a Gaussian distribution that defines its emission probabilities. As such, accepting the trade-off between generalization loss and efficiency gain, the `State` class directly embeds the parameters of the associated Gaussian distribution. This allows for faster retrieval and update of the parameters, and can be easily expanded to other probability distributions (for future work) by turning `State` into an abstract class and deriving distribution-specific `State` classes from it. It is worth noting that the emission probability is provided directly in logarithmic space, to avoid useless overheads and to speed up the computation.

3.2.4 Data generation

To generate data from the model, the framework Pomegranate is being used. It is a general HMM library, but the fact that the implementation language is Python makes any eventual speed comparison unfair, and thus it will not be used to apply standard algorithms to the model. Pomegranate can generate data in a simple and fast way; some Python scripts interface with it by defining a model that is coherent with the one used or estimated in WaHMM.

Specifically, `generate_data.py` reads the model from a file and generates an observation sequence of some length, optionally saving both the sequence and the generating state path to file in both plaintext and binary form; this allows to read the binary file for faster input processing, shrinking the running times considerably. At this time, k-means clustering can be performed on the data after it is generated.

3.2.5 Wavelet compression

The compression of the data happens through the `Compressor` class, which wraps an interface to HaMMLET. The observation sequence undergoes the Maxlet transform and is encapsulated in a `BreakpointArray`, a data structure that, given a certain threshold, subdivides the sequence into blocks. Each block gathers observations that are sufficiently close and thus are likely to have been generated while the model was in the same hidden state. A summary of each block statistics is stored in a parallel structure called `IntegralArray`. The combination of these two structures enables the creation of an interface for simple and efficient querying of consecutive blocks and their statistics summary. More information on how HaMMLET works can be found in [13].

The threshold used to define blocks can surely be subject to discussions: a low threshold will yield more blocks than needed, reducing the efficiency of compression; a high threshold will instead generate fewer bigger blocks that may group observations belonging to different states. The choice is taken from HaMMLET: the threshold is obtained by computing an estimate of the noise variance from the finest detail coefficients of the wavelet transform.

3.2.6 Logarithms summation

Some very useful declarations and functions are present in `utilities.hpp` and `commons.hpp`; other than functions to easily print and free matrices, the most important one is `sum_logarithms()`. Since the program operates in the logarithmic space, the probability products are converted into summations. A sum in the original space, though, has no simple logarithmic equivalent; this situation happens often, as in equation (2.4b).

Using symbols, the function needed to solve this problem is some function F so that, given two elements in logarithmic space, $\log(x)$ and $\log(y)$, it should produce:

$$\log(x + y) = F(\log(x), \log(y)) \quad (3.1)$$

The simplest option would be converting both elements back to the original space through exponentiation:

$$\log(x + y) = \log(e^{\log(x)} + e^{\log(y)}) \quad (3.2)$$

However, this solution can cause underflow when $\log(x)$ or $\log(y)$ are too negative. A simplification of this allows to write

$$\log(x + y) = \log(x) + \log(e^{\log(y) - \log(x)} + 1) \quad (3.3)$$

when $\log(x) > \log(y)$. Some workaround is required if any of the operands is $-\infty$, implying the presence of some `if` clauses before the actual computation. The `sum_logarithms()` function implements this operation using the `std::log1pf()` function from the standard library of C++ for a more efficient computation of the logarithm.

3.2.7 Saving results

The algorithms' results can be saved to file by adding the correct option when calling the program. Although the results directory may be changed, the default choice is the `/results` folder. After running the algorithms, WaHMM will write the results in properly named files. To produce an example, a set of standard results can also be generated using Pomegranate through the script `pomegranate_test.py`, more for a comparison with a different approach than for speed and performance.

3.2.8 Test automation

The process of testing the algorithms against generated data is of course central to the thesis work; given the huge number of tests to conduct, it makes sense to automate not only the testing process but also the extraction of meaningful results. The Python script `automated_test.py` takes care of the testing process through a very simple sequence of steps:

- generate the model according to the topology and number of states relative to the current case that is subject to test;
- generate data from the model and save to file both the sequence and the generating path;
- estimate a model for the training problem from the data and save it to file;
- execute the algorithms, time their execution and compute some performance measures;
- periodically save the test results to file for future analysis.

The Python script `results_aggregation.py` analyzes the results in an automated way, by not only computing the differences between the compressed and standard algorithms in the chosen metrics, but also producing their graphical representation in the form of plots, boxplots, and a merge of the two. The last one is the format used for the figures discussed in Chapter 4.

3.2.9 Other Python files

Several Python scripts have been written to ease WaHMM's usage and setup, not only with an interface that may be simpler to use but that also allows easier automation of the testing process. The `generate_model.py` script defines a model and saves it to file so that it can be used in other scripts and later be imported by WaHMM. `plot_data.py` and `plot_kmeans.py` are utilities to produce some outlook on the generated data and the estimated model respectively. `viterbi_comparison.py` simply compares Viterbi paths obtained with different algorithms to point out the differences, particularly useful when checking the accuracy of the compressed algorithms against the real generating path. `utilities_io.py` and `utilities_kmeans.py` simply provide helper functions for the other scripts.

3.3 Standard algorithms

This section will discuss the implementation of the standard algorithms; for an easier understanding, it can be useful to reference Figure 2.1 for a structural overview. All the three functions solving the problems accept some boolean flags to influence the type of output: **verbose** prints more information to the standard output; **silence** suppresses all the output messages; **tofile** specifies that the output results should be saved to file.

3.3.1 Evaluation problem

The evaluation problems can be solved by computing the forward matrix. Looking at Figure 2.1, each row is associated with a state; each column represents a time step. After the initialization described in (2.4a), the induction phase is constructed by initializing the forward variables to $-\infty$ for the current time step; then, the `sum_logarithms()` function is applied to accumulate the sum of the products in (2.4b); at last, the emission probability is added and the computation moves to the following time step. A simple sum of the forward variables over all states at the last time step yields the desired $P(O|\lambda)$ probability. Looking at the code, it is easy to confirm the computational complexity of $\mathcal{O}(T \cdot N^2)$ (as it was previously stated in Section 2.1.1).

3.3.2 Decoding problem

From a theoretical perspective, the decoding problem is very similar to the evaluation problem; in fact, the Viterbi algorithm differs from the forward algorithm in applying the *max* operator instead of the *sum*. The implementation is slightly more complex, requiring an additional matrix (named `statesViterbi`) with the same structure to hold the *argmax* results from (2.11b). After the initialization phase, a loop is used to join the computation with a classic maximum search on-the-fly; this allows to only iterate once over the states. The backtracking described in (2.11d) is then applied by appending each state of the path at the head of a list; in this way, a simple visit of the list yields the Viterbi state path. Again, the computational complexity of $\mathcal{O}(T \cdot N^2)$ that was discussed in Section 2.2 is confirmed by the implementation.

3.3.3 Training problem

The Baum-Welch algorithm is the most complex of the three. It is performed iteratively for a maximum amount of iterations or until the procedure is improving $P(O|\lambda)$ by an amount smaller than a predefined threshold (10^{-9} in the implementation). For this reason, to avoid a big overhead, most of the memory allocations happen in `training_problem_wrapper()`, which also handles the iterations of the algorithm. The result returned at each iteration is the evaluation probability $P(O|\lambda)$ relative to the model before the reestimation; this implementation choice avoids one useless computation of the forward matrix at each iteration of the Baum-Welch algorithm at the cost of performing one more iteration than needed after matching

the threshold.

The starting model supplied to the procedure can be directly read from a file using the option `estimate`; in the absence of a solid estimation from domain knowledge, it is possible to generate the model estimate through an automatic procedure that, given the number of states, selects some means and variances to be associated with states. The chosen method for this automatic procedure is the application of k-means clustering to the sequence; the proposed Python implementation uses the clusters' centroids as estimated means and computes the standard deviation with the canonical squared distance formula. The model is then saved to an output `data/kmeans_model` file, to be imported as a starting model when executing the training algorithm. This clustering operation can happen either directly when the data is generated, for faster execution, or later on, for externally supplied data.

The theory presented in Section 2.3 advocates for calculating every variable for all time steps: this implies a huge memory usage to store the forward and backward variables, plus γ_t and ε_t . The procedure does not need to store all the intermediate results: to save space, the backward variable and γ_t are calculated only for the considered time step; the other variables are simply summed to obtain the final cumulative sum over all time steps. Looking at the reestimation equations, it can be seen that the denominator $P(O|\lambda)$ gets simplified in all of them except for (2.17a); thus, it is more convenient to remove it from the equations and just add it to (2.17a) at the end.

Performing the reestimation of the Gaussian's mean using the formula in (2.21a) causes a problem; although working in logarithmic space eases the computations for the variable γ_t , it cannot be done when observations are negative; that is, if the reestimated mean is negative, it cannot be represented in the logarithmic space. To overcome this problem, the observation sequence must be rescaled by translating it by a value that is strictly greater than the minimum observed value so that only strictly positive observations are present. This rescaling is applied on-the-fly and for the mean reestimation only to avoid unnecessary computations; the added offset must be removed when moving out of the logarithmic space.

3.4 Compressed algorithms

This section contains a description of how the data compression changes the algorithms in the implementation. To ease the computations, the summary statistics of a block can be retrieved both individually (through `Compressor::blockSize()`, `::blockSum()` and `::blockSumSq()`) and together, through an *ad hoc* data structure called `blockdata` and the related function `Compressor::blockData()`.

A key part of the computation is calculating $e^{E_w(i)}$ for a block. As a first observation, working in logarithmic space makes it possible to avoid the exponentiation and just calculate $E_w(i)$. More interestingly, from (2.27a) it can be seen that $K(n_w, j)$ only depends on the state and the size of the block, and not on the actual observations. To make the algorithms faster, the `Model` class stores an array of hashmaps `Model::mKValues` of the $K(n_w, j)$ values per each state, adding entries as they are computed.

3.4.1 Evaluation problem

The forward algorithm is modified to deal with compressed data. The implementation is very similar to the uncompressed one: the `Compressor` class allows iterating through the blocks to perform the computation of the compressed forward variable following the equations (2.28a), (2.28b) and (2.28c).

3.4.2 Decoding problem

The compressed Viterbi algorithm iterates over the blocks in a similar fashion to the forward algorithm, as it happens for the uncompressed algorithms. To keep track of the actual Viterbi path, though, the block sizes must be recorded somewhere; in this way, during the backtracking phase, any block state can be associated with its size, specifying how long the sequence stays in that state.

3.4.3 Training problem

As shown in previous discussions, the compressed training algorithm is the one that differs the most from the uncompressed version. A first difference that is worth noting is that the backward matrix is computed fully using a compressed backward algorithm; the space overhead is more manageable in this case since the backward matrix has one column per block and not one per observation. The variables are computed according to the new definitions given in Section 2.5; the mean and variance computations are done outside of the logarithmic space, to avoid the problems caused by negative values of the sum that would require a translation of all values. An important consideration can be made for (2.27b): $K(n_w, j)$ can be precomputed at each iteration of the Baum-Welch algorithm. To find a balance between avoiding unnecessary computations and performing the same calculation every time, a map is used to keep track of the $K(n_w, j)$ values that have already been encountered; whenever a new one is found, the computation is made for all the states and stored in the map; this allows to speed up the compressed algorithm a little bit more, at the cost of some extra memory.

4

Results

This chapter presents the methodologies that have been chosen to evaluate the developed algorithms, discussing the parameters influence and how the results will be observed; this has the purpose to identify some conditions under which the compressed algorithms are worth using.

4.1 Testing setup

4.1.1 Choosing the parameters

To test the approach developed in this thesis, some HMMs need to be defined. Choo et al. [15] do a very good job defining some of the most frequent and useful model topologies in the field of reference, which is computational biology; nonetheless, they have a much broader scope and are found in many different applications. Three topologies will be analyzed:

- *fully connected model*, with every pair of states being connected and thus with the underlying graph being complete; the fully connected graph also includes the self-loops for each state;
- *circular model*, with an ergodic graph: the states are arranged in a circle, and a transition can only occur towards the same state or to the next one;
- *left-to-right model*, with an acyclic graph except for self-loops; the states are partially ordered and there are uniquely defined starting and ending states; transitions must be taken to visit the states following that order.

Each topology will be explored using different numbers of states, to see how the performance and accuracy varies. Specifically, models will be defined with 2, 3, and 5 states; this choice should give a perspective on how this parameter influences the results. About the observation sequence, both its length and an expected number of transitions should be discussed. For typical applications, the observations' length usually is in the order of 10^5 . Empirical experience from other research suggests that an adequate expected number of transitions is 10, which is a realistic magnitude for several kinds of sequences of interest. The state separations that will be tested are 10 different values of η as defined in (2.39), going from 0.1 to 1.0 in increases of 0.1. Testing algorithms against randomly generated data always presents the risk of introducing a non-deterministic bias into results. To contrast this problem, 100 sequences have been generated for each model. This allows not only to analyze some aggregate values to obtain some summarized information but also to study the spread of the performance measures and thus the results' stability. The chosen

aggregation method is the median since its robustness makes it possible to correctly ignore a solid number of outliers. The number of sequences has been chosen as a compromise between the robustness of the statistics and the feasibility of the execution times on a laptop.

4.1.2 Results evaluation

Before rushing to the results' discussion, it is useful to overview what type of information will be presented to analyze the results and the mathematical tools used to process them. Every measure defined in this section will be plotted against the state separation η .

Defining some measure of relative error is necessary, but it can be tricky, especially when facing zero values; the approach used in this thesis is to use the relative difference, defined as

$$d_r = \frac{|x - y|}{\max(|x|, |y|)}, \quad (4.1)$$

with the caveat of setting the error to zero when both x and y are zero.

For the states' estimation, the accuracy is evaluated using another measure borrowed from information theory: the Kullback-Leibler divergence (or KL-divergence), an indicator of how much a probability distribution differs from another one taken as a reference. A deeper explanation of this measure falls out of the scope of this thesis; the formula is

$$D_{KL}(P||Q) = \int_{-\infty}^{\infty} p(x) \log\left(\frac{p(x)}{q(x)}\right) dx \quad (4.2)$$

Now that the mathematical tools have been described, the different performance measures can be properly explained.

The evaluation problem produces the logarithm probability $P(O|\lambda)$ as defined in (2.4c); the performance of the compressed algorithm is measured as the relative difference from the standard result.

For the decoding problem, the indicator has been chosen to be the fraction of errors in the estimated generating path; thus, the performance measure is the relative difference between the compressed and the standard results.

The training problem is much more complex, and summarizing its results to obtain some indicator of how well the algorithm performs is not an easy task. The choice for this thesis is to compare both the compressed and the standard estimates against the real model, along three dimensions: the average accuracy on the states, measured with the KL-divergence; the average error on the transition log probabilities; the average error on the log probability of the initial distribution evaluated at the starting state. These three indicators are computed for both the compressed and the standard algorithms and then compared by simply computing their difference. Finally, the speedup will be analyzed with a simple ratio of the execution times, separately for the three algorithms; this measure includes the data input processing time to account for the overhead that the compressed algorithms require to elaborate the data.

To provide more insights about the plotted information, each plot has boxplots on the side to describe the distribution of the results. It is easy to see how discussing

the results will generate a lot of figures; for this reason, the main discussion will only concern the tests for the fully connected model; the other topologies will be discussed in relation with this model and their plots will be grouped at the end of this Chapter.

4.2 Evaluation problem

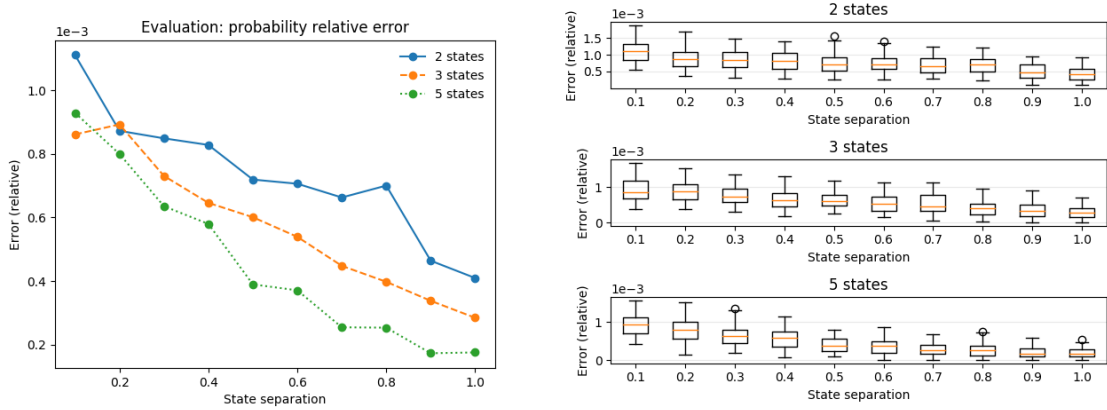


Figure 4.1: Relative difference between the $P(O|\lambda)$ log probabilities of the compressed and standard algorithms. Since the errors are on a big negative log probability, the actual error magnitude is approximately of 10^{-5000} .

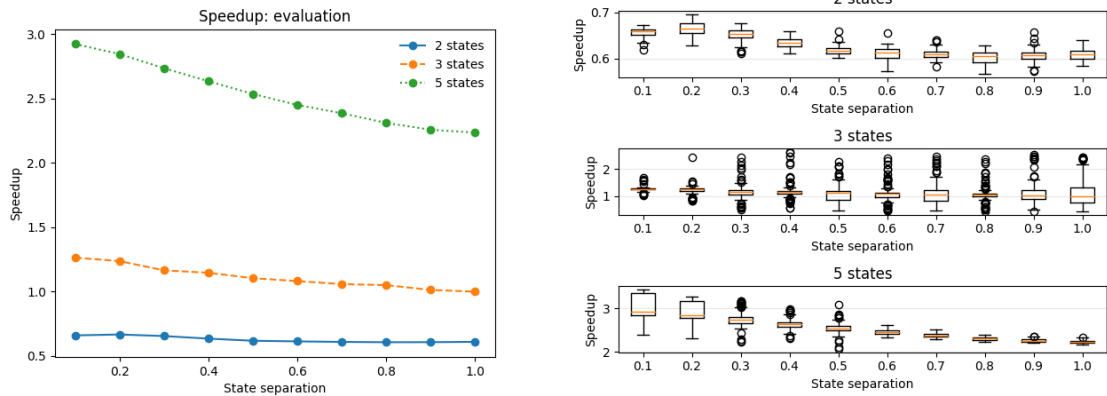


Figure 4.2: Speedup on the evaluation problem using the compressed algorithm, including the input data processing time.

The information plotted in Figure 4.1 is the relative difference of the logarithmic probability $P(O|\lambda)$ between the compressed and the standard algorithm. The first thing to notice is that as the state separation increases, the error decreases; this is expected since having more distinguishable states helps the compression process. In all cases, the error is relatively small and tightly spread, indicating that the compressed algorithm does a good job approximating the standard one on the evaluation

problem. Also, the error appears to decrease with a higher number of states in the model. To understand if the use of the compressed algorithm is worth it, Figure 4.2 shows the different speedups obtained; although they don't change much with the state separation, the main parameter affecting them is the number of states: the use of the compressed algorithm is thus advised only when the model has a relatively high number of states.

The circular and left-to-right topologies have slightly worse results, but they go through the same considerations made for the fully connected model.

4.3 Decoding problem

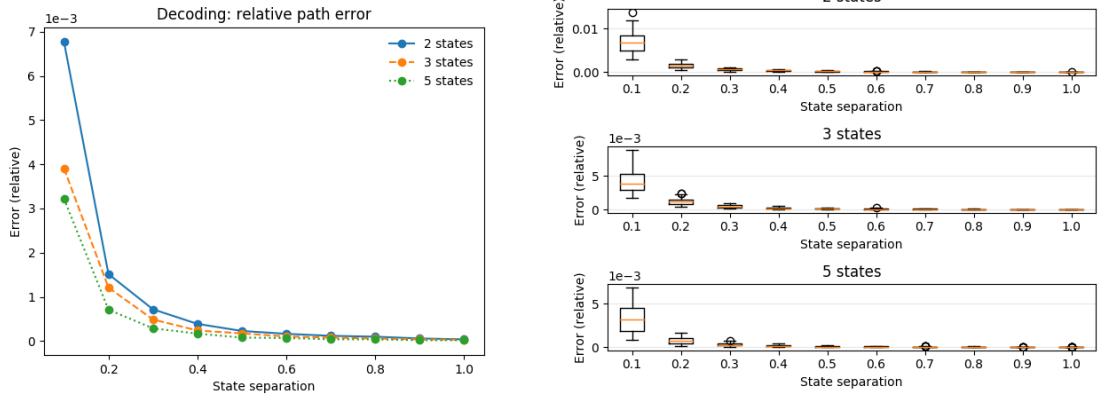


Figure 4.3: Relative difference between the fractions of errors in the estimated generating path.

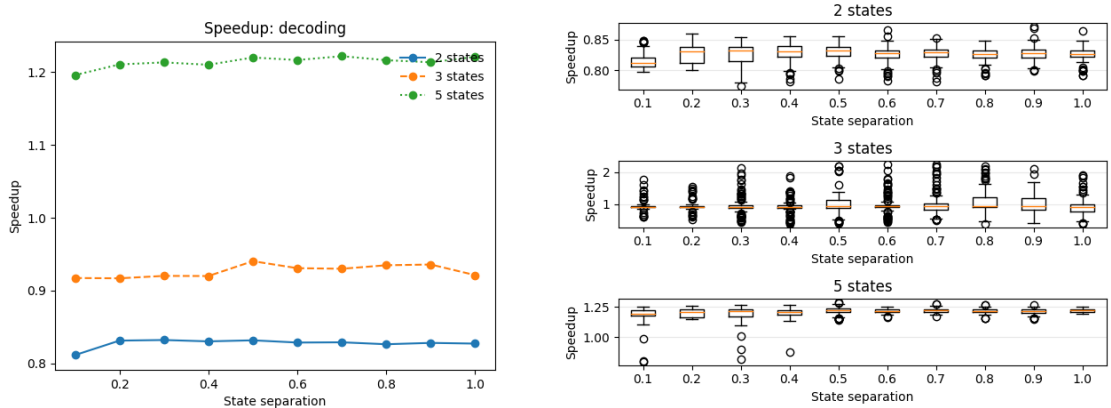


Figure 4.4: Speedup on the decoding problem using the compressed algorithm, including the input data processing time.

As properly explained in Section 2.2, the decoding problem is quite similar to the evaluation problem with the forward algorithm being partly modified. Thus, the expected results should more or less align to the performances in the evaluation

problem. As Figure 4.3 shows, the general trend is the same. When the state separation is extremely low, the uncertainty on the values is high; thus the actual ordering between the states, although consistent with Figure 4.1, may vary slightly in certain points. Although the performance indicator seems good, Figure 4.4 shows how the compressed algorithm is generally less efficient than the standard one when the model has a small number of states; moreover, the data has a very high spread, suggesting that the actual speedup depends a lot on the data being generated.

The model that scores best in the decoding problem is the left-to-right model, which intuitively makes sense since each state only has one allowed transition; again, though, the results are pretty similar and thus the compressed algorithm should only be used when the number of states is high enough to achieve an actual speedup.

4.4 Training problem

The training problem is the most complex of the three, as it has been said many times at this point. The first plot of interest is in Figure 4.5, showing how well the states are estimated by the compressed algorithm compared to the standard one. While the overall trend is that of a constant small difference for most values of state separation, it is noteworthy that the plotted values are negative when the states are not well-separated. This means that the compressed algorithm is more accurate than the standard one, and this accuracy appears to increase with the number of states if the state separation is small enough; the main problem is that, for an increasing number of the states, the spread of the results starts becoming very high. Despite this, for a 5-states model that is not enough to cause problems or big inaccuracies.

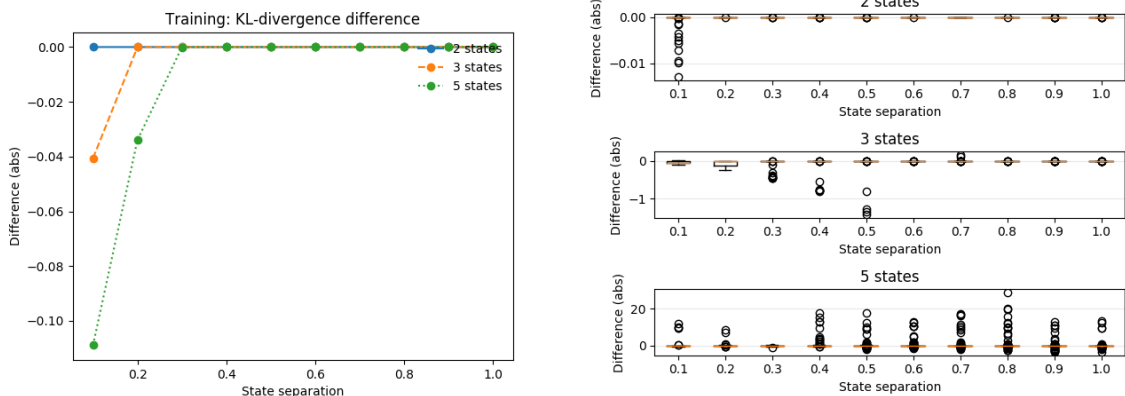


Figure 4.5: Difference between the average KL-divergence for the compressed and standard algorithms.

Similar considerations can be made for both the transition probabilities and the initial distribution estimations; the compressed training does an overall better job by a small margin, that gets more consistent when the state separation is very low.

4. Results

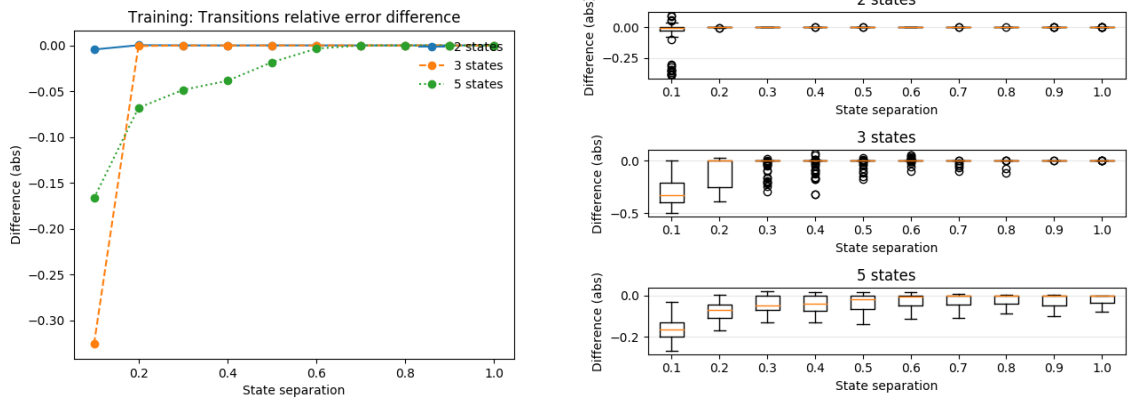


Figure 4.6: Relative difference between the average error on the log transition probabilities of the compressed and standard algorithms.

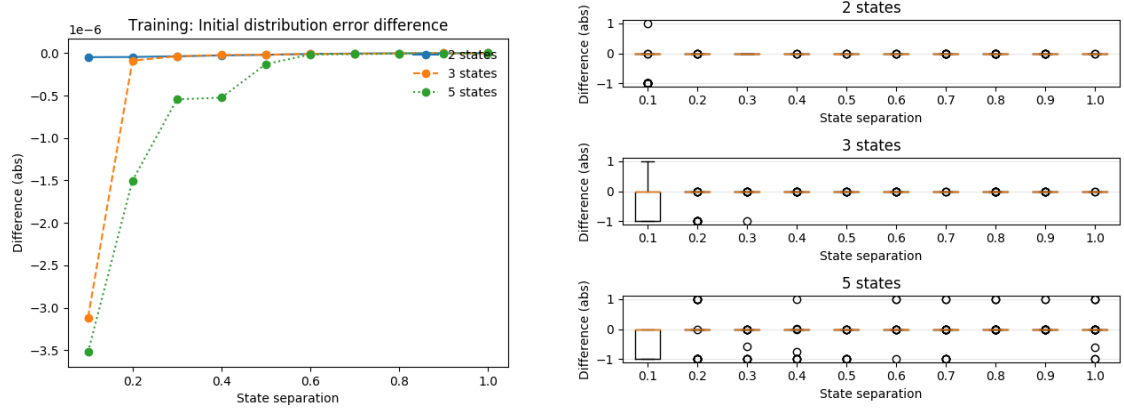


Figure 4.7: Relative difference between the average error on the log probability of the starting state for the compressed and standard algorithms.

The results for other topologies resemble the presented ones quite well. The generally high instability of the results can make the appearance of the plot less meaningful; but when taking a deeper look at the data, the distribution of the results is mostly skewed towards favoring the compressed algorithm over the standard one. Looking at the speedup in Figure 4.8, there is a very noticeable performance gain that is higher when the state separation is low and the number of states is high. In general, the speedup tends to be extremely high, giving a solid reason to use the compressed training algorithm.

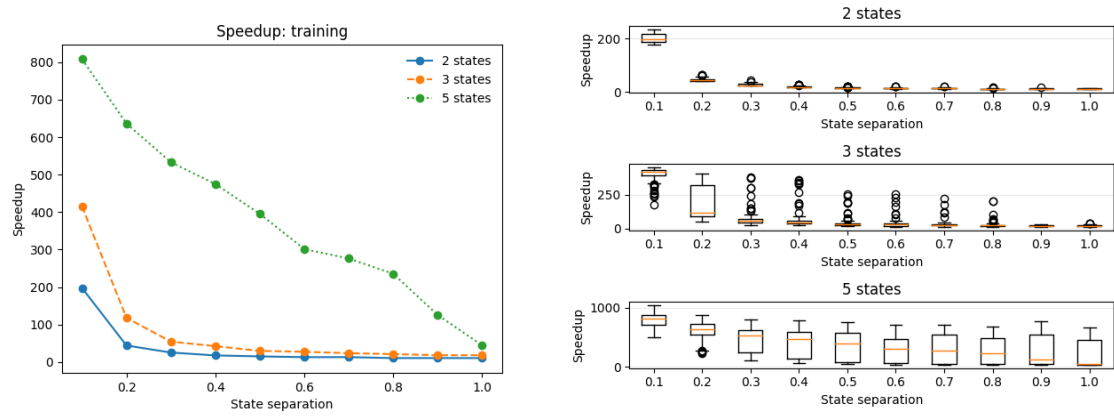


Figure 4.8: Speedup on the training problem using the compressed algorithm, including the input data processing time.

4.5 Circular topology

This section presents the results of the testing process performed on the *circular* model topology. The compressed results are very close to the ones produced by the standard algorithms; comparing the speedups in Figure 4.10 or Figure 4.12 with the respective plots for the fully connected model, it is worth noting that the compressed evaluation and decoding algorithms yield a smaller speedup. This is not true for the training algorithm, that maintains the considerable speedup that was achieved for the fully connected model.

4.5.1 Evaluation problem

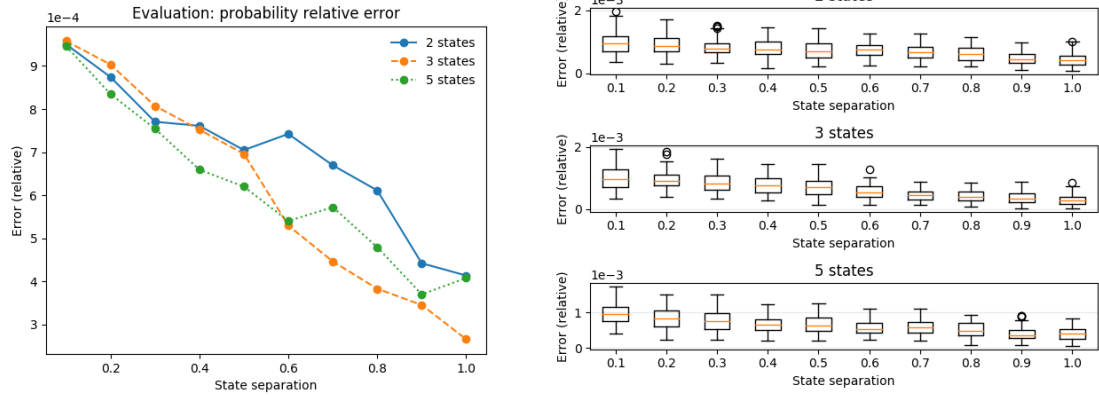


Figure 4.9: Relative difference between the $P(O|\lambda)$ log probabilities of the compressed and standard algorithms.

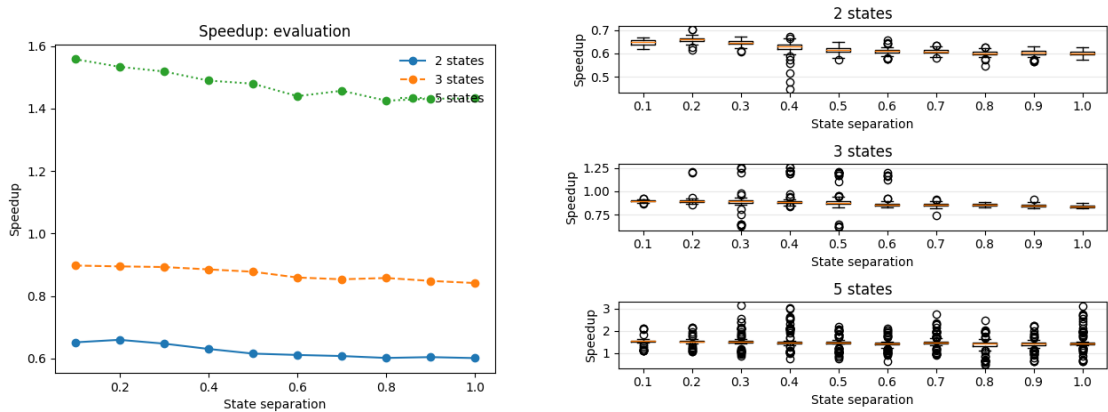


Figure 4.10: Speedup on the evaluation problem using the compressed algorithm, including the input data processing time.

4.5.2 Decoding problem

In Figure 4.12 the speedup for an increasing number of states is shown to be less and less stable; however, the outliers in the boxplot are skewed mostly towards higher speedups. This number of outliers suggests that having more tests could yield better information on the actual results distribution.

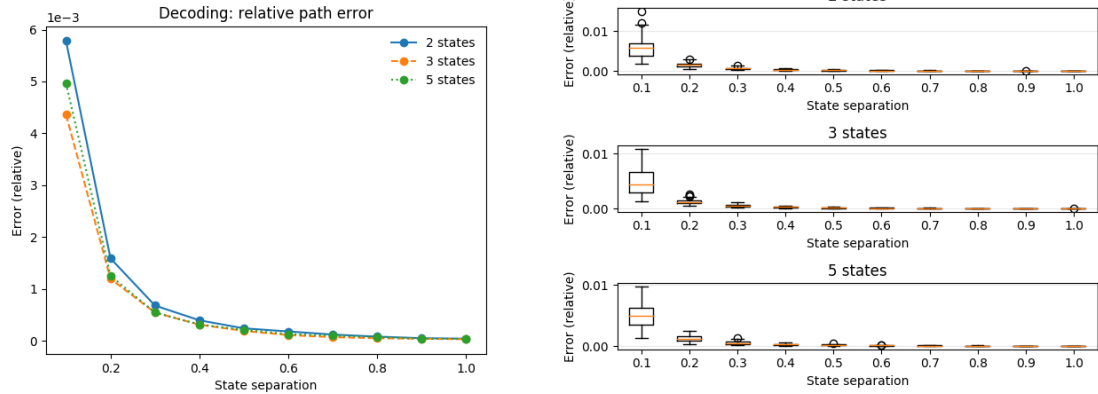


Figure 4.11: Relative difference between the fractions of errors in the estimated generating path.

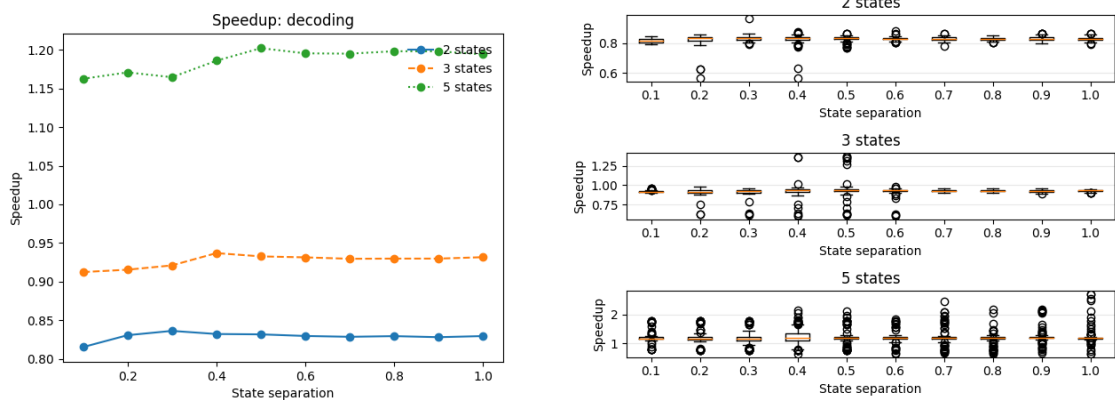


Figure 4.12: Speedup on the decoding problem using the compressed algorithm, including the input data processing time.

4.5.3 Training problem

The training results are very similar to the plots for the fully connected model, thus they follow a similar discussion.

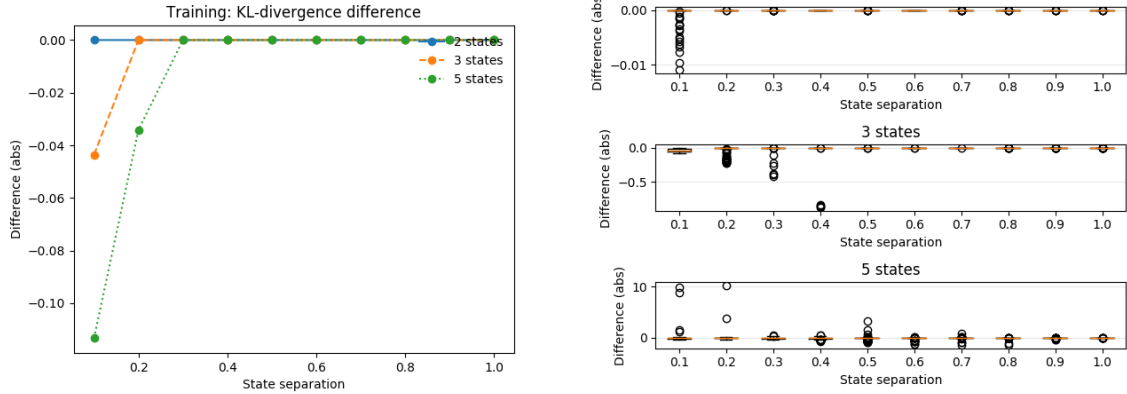


Figure 4.13: Difference between the average KL-divergence for the compressed and standard algorithms.

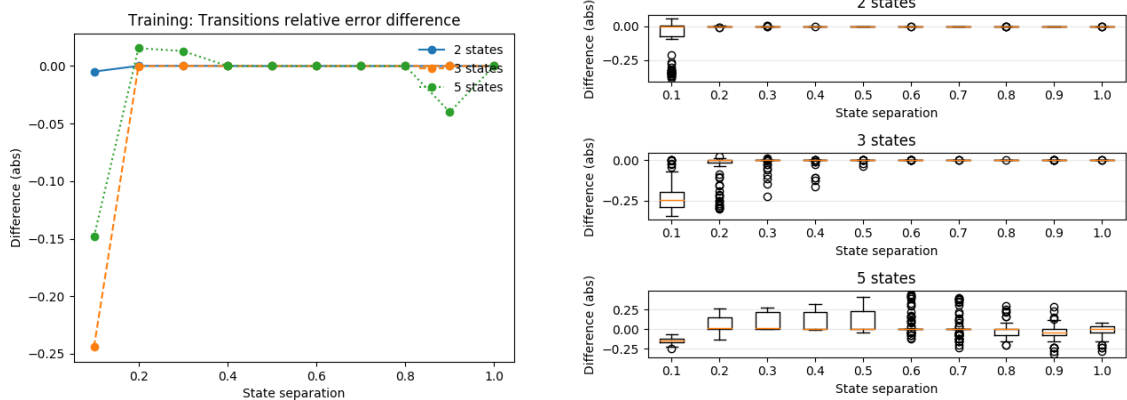


Figure 4.14: Relative difference between the average error on the log transition probabilities of the compressed and standard algorithms.

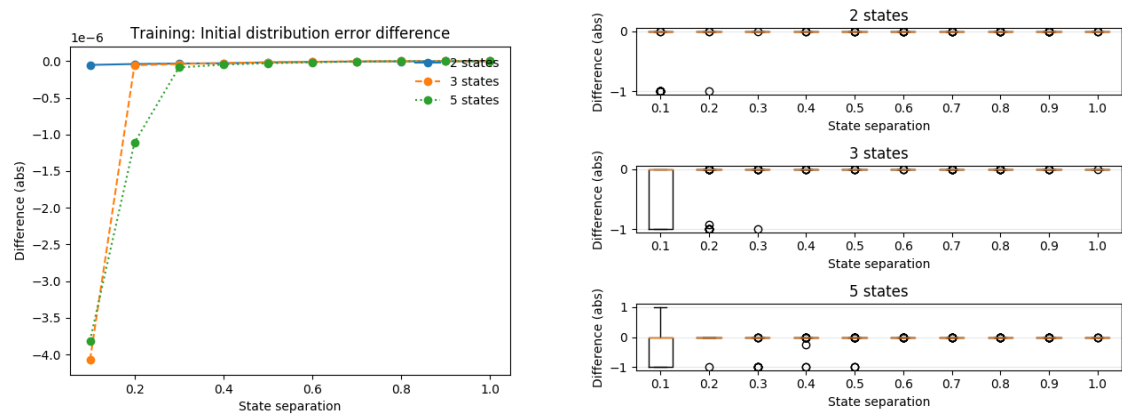


Figure 4.15: Relative difference between the average error on the log probability of the starting state for the compressed and standard algorithms.

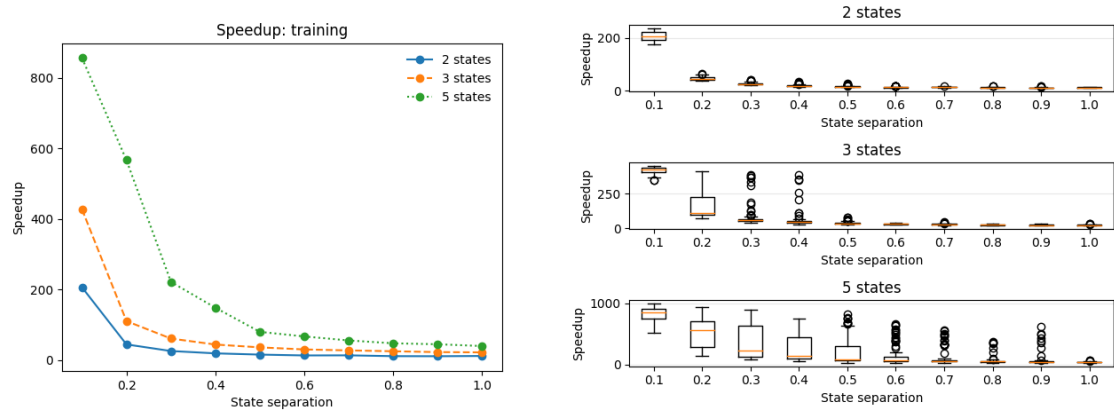


Figure 4.16: Speedup on the training problem using the compressed algorithm, including the input data processing time.

4.6 Left-to-right topology

This section presents the results of the testing process performed on the *left-to-right* model topology.

4.6.1 Evaluation problem

The evaluation plotted results are more separated than in the previous cases; it is important to remember that the scale of the error is generally very small, thus this difference is not extremely relevant.

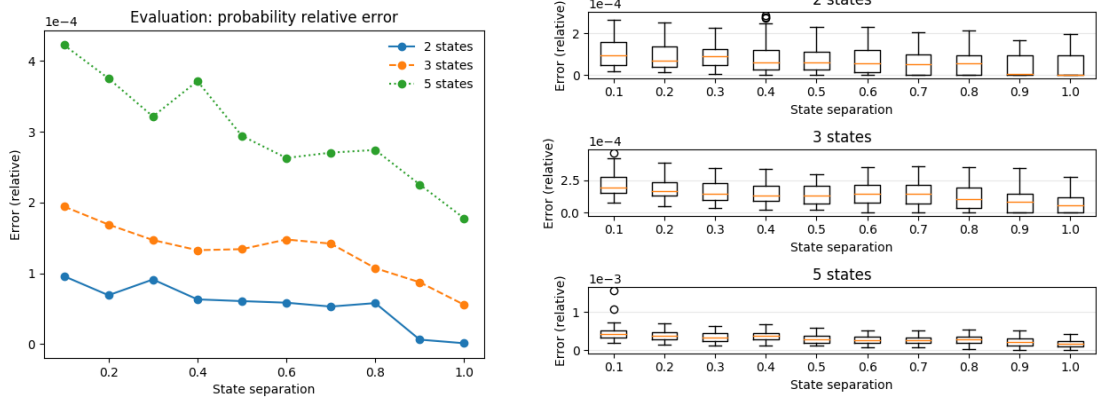


Figure 4.17: Relative difference between the $P(O|\lambda)$ log probabilities of the compressed and standard algorithms.

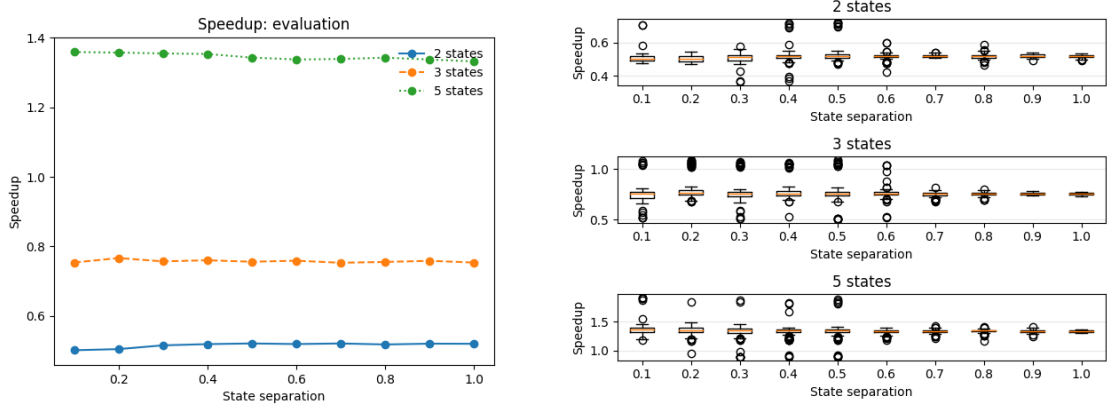


Figure 4.18: Speedup on the evaluation problem using the compressed algorithm, including the input data processing time.

4.6.2 Decoding problem

The decoding algorithm for this topology compared to the other ones is more spread out for a lower number of states and less spread out for a higher number of states. In the end, the plots in Figure 4.20 appear a bit more stable; the key point is still that the speedup becomes greater than 1 for a number of states greater than 5.

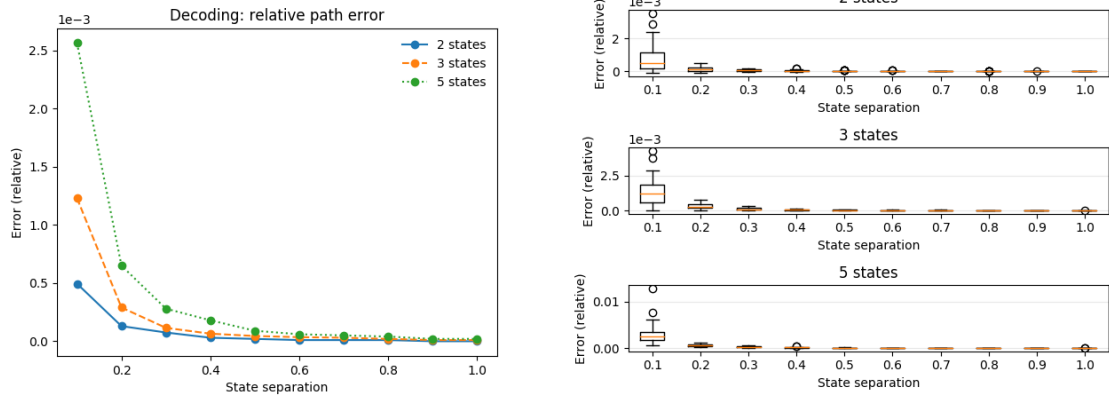


Figure 4.19: Relative difference between the fractions of errors in the estimated generating path.

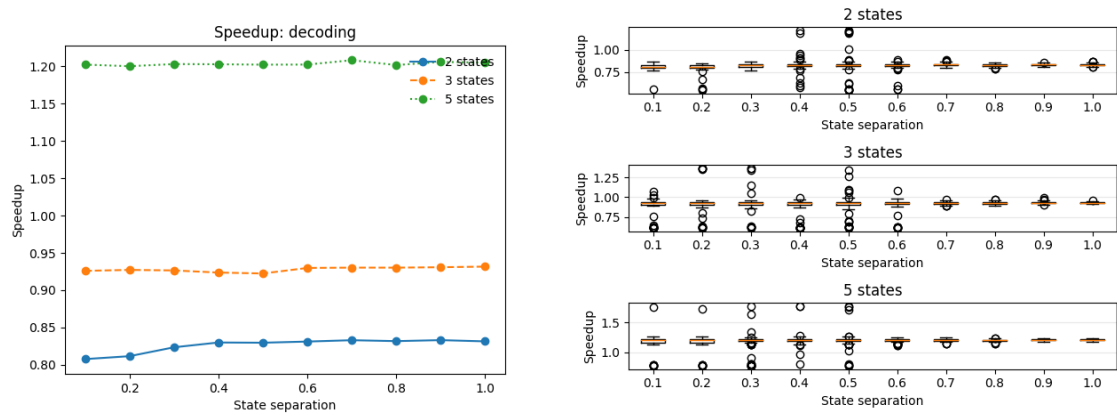


Figure 4.20: Speedup on the decoding problem using the compressed algorithm, including the input data processing time.

4.6.3 Training problem

These results are visually somewhat different from the previous ones. It is important to always consider the scale at which the error is represented and the number of outliers present in the data; in fact, a more attentive look shows that the results for the training algorithm with the left-to-right model topology are not that different from the other ones. It is worth noting that the speedups are generally slightly higher, at least for an increasing amount of states.

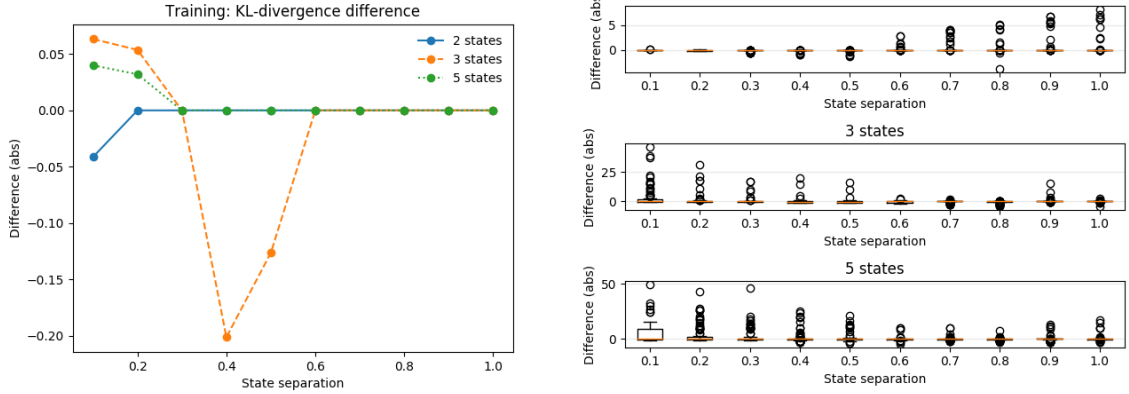


Figure 4.21: Difference between the average KL-divergence for the compressed and standard algorithms.

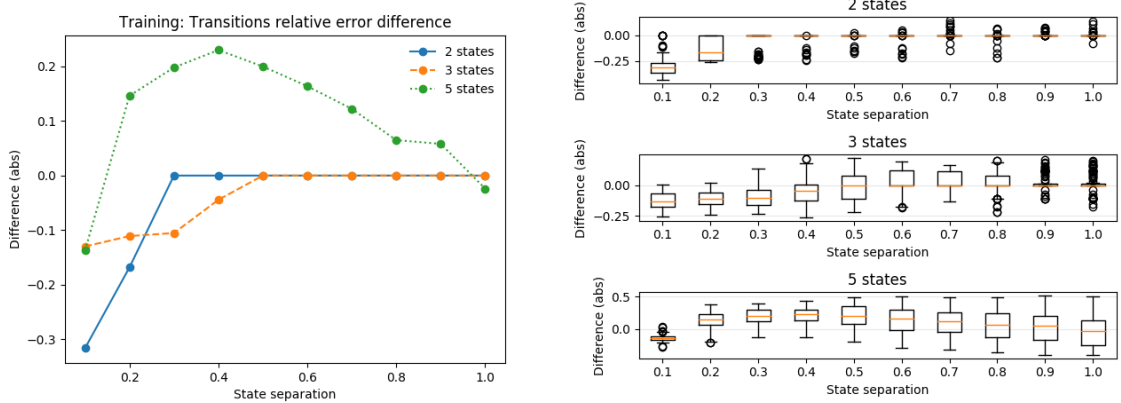


Figure 4.22: Relative difference between the average error on the log transition probabilities of the compressed and standard algorithms.

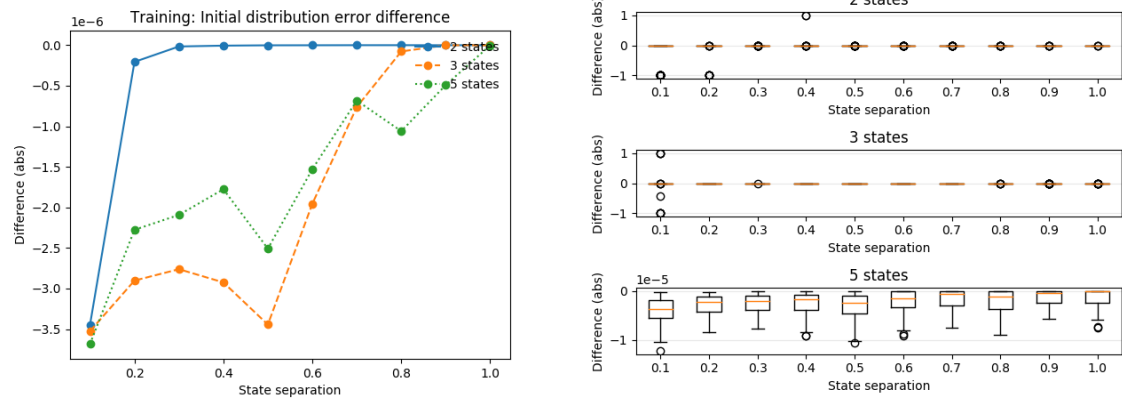


Figure 4.23: Relative difference between the average error on the log probability of the starting state for the compressed and standard algorithms.

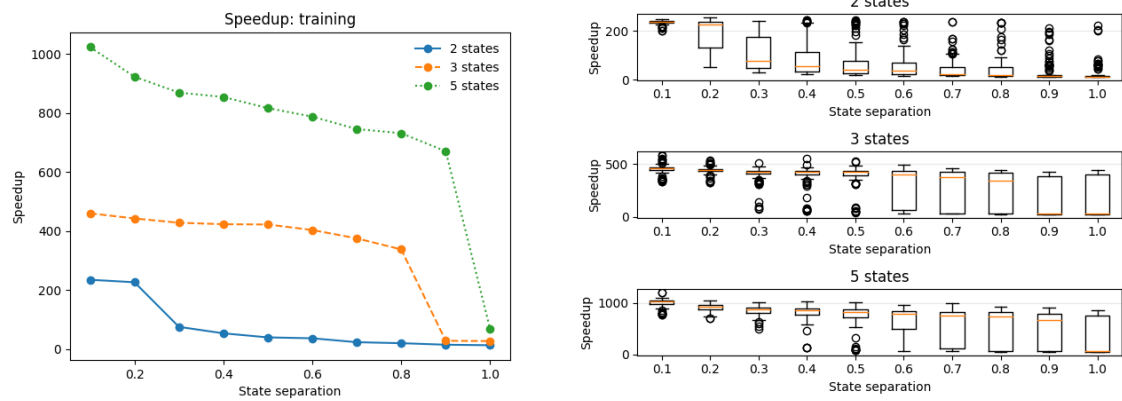


Figure 4.24: Speedup on the training problem using the compressed algorithm, including the input data processing time.

5

Conclusion

5.1 Main takeaways

The previous chapter presented a description of the results on a case-by-case basis; after that, it is useful to summarize the insights that have been extracted from the testing process.

The first thing to point out is that the results are pretty consistent over all the tested topologies; this makes it possible to draw some general conclusions that are independent of the topology of the model. Overall, the compressed algorithms perform with extremely high precision. Specifically, the evaluation and the decoding compressed algorithms perform as well as the standard ones regardless of the states' separation. They are slower than their standard counterpart for models with a low number of states, because of the overhead caused by the data compression; fortunately, the speedup gets better for models with a higher number of states, when the running times would start to grow more and more. The training algorithm is the one that performs best: it performs better than the standard algorithm for a low state separation, and the achieved speedup is extremely high.

It is important to remember that there are elements that have been assumed constant throughout the testing process, such as the self-transition probabilities or the sequence length. In the context of big data, for sequences that are much longer than the ones used in testing, the compressed algorithms certainly perform even better than what is shown in Chapter 4.

5.2 Wrapping up

WaHMM gives the opportunity to apply either standard or compressed algorithms to solve the evaluation, decoding, and training problems with an efficient C++ implementation. In the context of scientific research, this thesis will hopefully serve as another confirmation that wavelet compression can work really well to allow hidden Markov models to scale to big sequence data. In particular, it enables the training of a hidden Markov model on a commodity laptop instead of requiring more complex machinery, since the training might be even thousands of times faster.

Bibliography

- [1] I. P. Waldmann, “On signals faint and sparse: The ACICA algorithm for blind de-trending of exoplanetary transits with low signal-to-noise,” *Astrophys. J.*, vol. 780, p. 23, 2014.
- [2] K. R. e. a. Wiedenhoeft J., Cagan A., “Bayesian localization of CNV candidates in WGS data within minutes,” *Algorithms Mol Biol*, vol. 14, no. 20, 2019.
- [3] I. R. Sipos, A. Ceffer, and J. Levendovszky, “Parallel optimization of sparse portfolios with ar-hmms,” *Computational Economics*, vol. 49, pp. 563–578, Apr 2017.
- [4] S. Y. M. Gales, “The Application of Hidden Markov Models in Speech Recognition,” *Foundations and Trends in Signal Processing*, vol. 1, no. 3, pp. 195–304, 2007.
- [5] Z.-U. M. Mozes S., Weimann O., “Speeding Up HMM Decoding and Training by Exploiting Sequence Repetitions,” in *Combinatorial Pattern Matching* (B. Ma and K. Zhang, eds.), (Berlin, Heidelberg), Springer Berlin Heidelberg, 2007.
- [6] M. P. Mahmud and A. Schliep, “Speeding Up Bayesian HMM by the Four Russians Method,” in *Algorithms in Bioinformatics* (T. M. Przytycka and M.-F. Sagot, eds.), (Berlin, Heidelberg), pp. 188–200, Springer Berlin Heidelberg, 2011.
- [7] S. A. Wiedenhoeft J., Brugel E., “Fast Bayesian Inference of Copy Number Variants using Hidden Markov Models with Wavelet Compression,” *PLoS Comput Biol*, vol. 12, no. 5, 2016.
- [8] L. R. Rabiner, “A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition,” *Proc. IEEE*, vol. 77, no. 2, pp. 261–266, 1989.
- [9] A. B. Romeo, C. Horellou, and J. Bergh, “A wavelet add-on code for new-generation N-body simulations and data de-noising (JOFILUREN),” *Monthly Notices of the Royal Astronomical Society*, vol. 354, pp. 1208–1222, 11 2004.
- [10] L. Liporace, “Maximum likelihood estimation for multivariate observations of Markov sources,” *IEEE Transactions on Information Theory*, vol. 28, no. 5, pp. 729–734, September 1982.
- [11] B. H. Juang, “Maximum-likelihood estimation for mixture multivariate stochastic observations of Markov chains,” *AT&T Technical Journal*, vol. 64, no. 6, pp. 1235–1249, July-Aug. 1985.
- [12] B. Juang, S. Levinson, and M. Sondhi, “Maximum likelihood estimation for multivariate mixture observations of markov chains,” *IEEE Transactions on Information Theory*, vol. 32, pp. 307–309, 3 1986.
- [13] J. Wiedenhoeft, *Dynamically compressed Bayesian hidden Markov models using Haar wavelets*. PhD thesis, Rutgers, The State University of New Jersey, 2018.

- [14] “Exact real arithmetic - haskellwiki - https://wiki.haskell.org/exact_real_arithmetic.”
- [15] . Z. L. Choo K. H., Tong J. C., “Recent applications of hidden markov models in computational biology,” *Genomics, Proteomics & Bioinformatics*, vol. 2, 5 2004.