

Spectral Learning of Mixture of Hidden Markov Models

Authored by:

Paris Smaragdis
Cem Subakan
Johannes Traa

Abstract

In this paper, we propose a learning approach for the Mixture of Hidden Markov Models (MHMM) based on the Method of Moments (MoM). Computational advantages of MoM make MHMM learning amenable for large data sets. It is not possible to directly learn an MHMM with existing learning approaches, mainly due to a permutation ambiguity in the estimation process. We show that it is possible to resolve this ambiguity using the spectral properties of a global transition matrix even in the presence of estimation noise. We demonstrate the validity of our approach on synthetic and real data.

1 Paper Body

Method of Moments (MoM) based algorithms [1, 2, 3] for learning latent variable models have recently become popular in the machine learning community. They provide uniqueness guarantees in parameter estimation and are a computationally lighter alternative compared to more traditional maximum likelihood approaches. The main reason behind the computational advantage is that once the moment expressions are acquired, the rest of the learning work amounts to factorizing a moment matrix whose size is independent of the number of data items. However, it is unclear how to use these algorithms for more complicated models such as Mixture of Hidden Markov Models (MHMM). MHMM [4] is a useful model for clustering sequences, and has various applications [5, 6, 7]. The E-step of the Expectation Maximization (EM) algorithm for an MHMM requires running forwardbackward message passing along the latent state chain for each sequence in the dataset in every EM iteration. For this reason, if the number of sequences in the dataset is large, EM can be computationally prohibitive. In this paper, we propose a learning algorithm based on the method of moments for MHMM. We use the fact that an MHMM can be expressed as an HMM with block diagonal transition matrix. Having made that observation,

we use an existing MoM algorithm to learn the parameters up to a permutation ambiguity. However, this doesn't recover the parameters of the individual HMMs. We exploit the spectral properties of the global transition matrix to estimate a de-permutation mapping that enables us to recover the parameters of the individual HMMs. We also specify a method that can recover the number of HMMs under several spectral conditions.

2 2.1

Model Definitions Hidden Markov Model

In a Hidden Markov Model (HMM), an observed sequence $\mathbf{x} = \mathbf{x}_{1:T} = \{x_1, \dots, x_t, \dots, x_T\}$ with $x_t \in \mathcal{X}$ is generated conditioned on a latent Markov chain $\mathbf{r} = \mathbf{r}_{1:T} = \{r_1, \dots, r_t, \dots, r_T\}$, with $r_t \in \{1, \dots, M\}$.

The HMM is parameterized by an emission matrix $\mathbf{O} \in \mathbb{R}^{|\mathcal{X}| \times M}$, a transition matrix $\mathbf{A} \in \mathbb{R}^{M \times M}$ and an initial state distribution $\pi \in \mathbb{R}^M$. Given the model parameters $\theta = (\mathbf{O}, \mathbf{A}, \pi)$, the likelihood of an observation sequence $\mathbf{x}_{1:T}$ is defined as follows: $p(\mathbf{x}_{1:T} | \theta) =$

$$\begin{aligned} & \sum_{\mathbf{r}_{1:T}} p(\mathbf{x}_{1:T}, \mathbf{r}_{1:T} | \theta) = \\ & \sum_{\mathbf{r}_{1:T}} \prod_{t=1}^T p(x_t | r_t, \mathbf{O}) p(r_t | r_{t-1}, \mathbf{A}) \\ & = \sum_{\mathbf{r}_{1:T}} \prod_{t=1}^T \text{diag}(p(\mathbf{x}_T | \mathbf{r}_T, \mathbf{O})) \mathbf{A} \text{diag}(p(\mathbf{x}_1 | \mathbf{r}_1, \mathbf{O})) \pi = \\ & \sum_{\mathbf{r}_{1:T}} \prod_{t=1}^T \text{diag}(p(\mathbf{x}_t | \mathbf{r}_t, \mathbf{O})) \mathbf{A} \text{diag}(p(\mathbf{x}_1 | \mathbf{r}_1, \mathbf{O})) \pi \end{aligned}$$

where $\mathbf{1}_M \in \mathbb{R}^M$ is a column vector of ones, we have switched from index notation to matrix notation in the second line such that summations are embedded in matrix multiplications, and we use the MATLAB colon notation to pick a row/column of a matrix. Note that $\mathbf{O}(x_t) := p(x_t | :, \mathbf{O})$. The model parameters are defined as follows: $\pi(u) = p(r_1 = u | r_0) = p(r_1 = u)$

initial latent state distribution

$$\mathbf{A}(u, v) = p(r_t = u | r_{t-1} = v), t \geq 2$$

latent state transition matrix

$$\mathbf{O}(:, u) = E[x_t | r_t = u]$$

emission matrix

The choice of the observation model $p(x_t | r_t)$ determines what the columns of \mathbf{O} correspond to: \bullet Gaussian: $p(x_t | r_t = u) = \mathcal{N}(x_t; \mu_u, \sigma^2)$

\bullet

$$\mathbf{O}(:, u) = E[x_t | r_t = u] = \mu_u$$

$$\bullet \text{ Poisson: } p(x_t | r_t = u) = \text{PO}(x_t; \mu_u)$$

\bullet

$$\mathbf{O}(:, u) = E[x_t | r_t = u] = \mu_u$$

$$\bullet \text{ Multinomial: } p(x_t | r_t = u) = \text{Mult}(x_t; \mu_u, S)$$

\bullet

3

Traditionally, the parameters of an MHMM are learned with the Expectation-maximization (EM) algorithm. One drawback of EM is that it requires a good initialization. Another issue is its computational requirements. In every iteration, one has to perform forward-backward message passing for every sequence, resulting in a computationally expensive process, especially when dealing with large datasets. The proposed MoM approach avoids the issues associated with EM by leveraging the information in various moments computed from the data. Even these moments, which can be computed efficiently, the computation time of the learning algorithm is independent of the amount of data (number of sequences and their lengths). Our approach is mainly based on the observation that an MHMM can be seen as a single HMM with a block-diagonal transition matrix. We will first establish this proposition and discuss its implications. Then, we will describe the proposed learning algorithm. 3.1

[illegible]

?K ?K

(4)

$$k=1$$

Tn Y

 $t=1$

We see that the state space of an MHMM consists of K disconnected regimes. For each sequence sampled from the MHMM, the first latent state r_1 determines what region the entire latent state sequence lies in. 3.2

In the previous section, we showed the equivalence between the MHMM and an HMM with a blockdiagonal transition matrix. Therefore, it should be possible to use an HMM learning algorithm such as spectral learning for HMMs [1, 2] to find the parameters of an MHMM. However, the true global parameters θ are recovered inexactly due to noise : $\hat{\theta} \neq \theta$ and state indexing ambiguity via a P, A^*P permutation mapping $P: \mathcal{S} \rightarrow \mathcal{S}$. Consequently,

corresponding to the eigenvalues which are 1. This algorithm corresponds to the noiseless case $A \rightarrow P$. In practice, the output of the learning algorithm \hat{P} is $A \rightarrow P$ and the clear structure in Equation (6) no longer holds in (A) , as $e \rightarrow ?$, as illustrated in the bottom row of Figure 1. We can see that the three-cluster structure no longer holds for large e . Instead, the columns of the transition matrix converge to a global stationary distribution.

3.2.2 Estimating the permutation in the presence of noise

In the general case with noise, we lose the spectral property that the global transition matrix has K eigenvalues which are 1. Consequently, the algorithm described in Section 3.2.1 cannot be

e: 1
e: 5
e: 10
e: 20
e: 1
rt+1
e: 5
e: 10
e: 20
rt
rt
rt
rt+1 rt
rt
rt
rt
rt
e: 1
e: 5
e: 10
e: 20
rt
rt
rt
rt
rt+1

Figure 1: (Top left) Block-diagonal transition matrix after e -fold exponentiation. Each block converge to its own stationary distribution. (Top right) Same as above with permutation. (Bottom) Corrupted and permuted transition matrix after exponentiation. The true number $K = 3$ of HMMs is clear for intermediate values of e , but as $e \rightarrow ?$, the columns of the matrix converge to a global stationary distribution. applied directly to make $A \rightarrow P$ block diagonal. In practice, the estimated transition matrix has only one eigenvalue with unit magnitude and $\lim_{e \rightarrow \infty} (A \rightarrow P)$ converges to a global stationary distribution. However, if the noise is sufficiently small, a depermutation mapping P clusters K can be successfully estimated. We now specify the

Eigenvalue Index

(Right) Spectral

$$\begin{aligned} & \max_{k \in \{1, \dots, K\}} \\ & \text{---} L_{2,k} \text{ --- and } \arg \max_{K=0} \\ & e_0 \text{ ---} ? \quad K e_0 \text{ ---} ? \quad e_0 \text{ ---} ? \quad K + 1 \quad K ? 1 \\ & = K, \text{ for } K = 0 ? \end{aligned}$$

we have $\arg \max_{K \in \mathcal{K}} K = K$.

Proposed Algorithm

8

4.4.1

Experiments Effect of noise on depermutation algorithm

We have tested the algorithm's performance with respect to amount of data. We used the parameters $K = 3$, $M = 4$, $L = 20$, and we have 2 sequences with length T for each cluster. We used a Gaussian observation model with unit observation variance and the columns of the emission matrices $O_{1:K}$ were drawn from zero mean spherical Gaussian with variance 2. Results for 10 uniformly 6

Algorithm 1 Spectral Learning for Mixture of Hidden Markov Models Inputs: $x_{1:N}$: Sequences, M, K : total number of states of global HMM. b, A, b, A : MHMM parameters Output: $\hat{b} = O_{1:K}^{-1} O_{1:K}$ Method of Moments Parameter Estimation $\hat{P}, \hat{A}^*P(O) = \text{HMM Method of Moments}(x_{1:N}, M, K)$ Depermutation Find eigenvalues of \hat{A}^*P

Exponentiate eigenvalues for each discrete value e in a sufficiently large range. b as the eigenvalue with largest longevity. Identify K, b reconstruction \hat{A}_r via eigendecomposition. Compute rank- K, b clusters to find a depermutation mapping P, e via cluster labels. Cluster the columns of \hat{A}_r with K, P, P^* and \hat{A}^* according to P, e Depermute O, b^*P and \hat{A}^*P Form \hat{P} by choosing corresponding blocks from depermuted O, b Return \hat{P} . Euclidean Distance vs Sequence Length Euc. Dist.

2
1
0
10
120
230
340
450
560
670
3
3
780
890
1000
T
3
3
3
3
3
3
3
3
3

Figure 3: Top row: Euclidean distance vs T . Second row: Noisy input matrix. Third row: Noisy reconstruction \hat{A}_r . Bottom row: Depermuted matrix,

numbers at the bottom indicate the estimated number of clusters. spaced sequence lengths from 10 to 1000 are shown in Figure 3. On the top row, we plot the total error (from centroid to point) obtained after fitting k-means with true number of HMM clusters. We can see that the correct number of clusters $K = 3$ as well as the block-diagonal structure of the transition matrix is correctly recovered even in the case where $T = 20$. 4.2

Amount of data vs accuracy and speed

We have compared clustering accuracies of EM and our approach on data sampled from a Gaussian emission MHMM. Means of each state of each cluster is drawn from a zero mean unit variance Gaussian, and observation covariance is spherical with variance 2. We set $L = 20$, $K = 5$, $M = 3$. We used uniform mixing proportions and uniform initial state distribution. We evaluated the clustering accuracies for 10 uniformly spaced sequence lengths (every sequence has the same length) between 20 and 200, and 10 uniformly spaced number of sequences between 1 and 100 for each cluster. The results are shown in Figure 4. Although EM seems to provide higher accuracy on 7

Accuracy (%) of EM algorithm
7 10 13 15 17 20 22 25
60 100 88 81 100100100100100100
60 80 100100100100 80 80 100100
3
5
7
9 12 14 16 18 20 23
78 80 95 100 98 100100100100 79 100
78 60 80 80 100100 80 100100100100
78 1
4
6
8 11 13 14 16 18 20
20 62 86 100100100 78 100100 80
80 100100100100100100100 80 100
2
4
6
7
9 11 13 14 16 18
56 80 77 82 81 60 100 88 100100 77
56 60 100100100100100100100100100
56 1
4
5
6
8 10 11 13 14 16
80 100 71 100100100100100100100
1

3
 5
 6
 7
 8
 9 11 12 13
 34 80 83 66 79 97 69 100 80 78 100
 34 40 100100100100100100100100100
 34 1
 3
 4
 5
 6
 7
 8
 9 10 11
 80 82 97 65 61 69 69 88 82 80
 60 100100100100100100100100100100
 1
 3
 3
 4
 5
 5
 6
 7
 7
 8
 12 20 65 73 76 78 77 77 78 63 86
 12 80 100100100100100100100100100
 12 2
 3
 3
 3
 3
 4
 4
 5
 5
 6
 1 20 53 68 58 73 79 76 88 58 78 10 31 73 116 158 200 T
 1 60 100100100100100100100100100 10 31 73 116 158 200 T
 1 2 2 2 3 3 3 3 3 3 10 31 73 116 158 200 T
 100
 75
 614 1616 2433 2423 2404 3332 5849 4915 6890

56
 573 1056 1418 3074 2030 3603 5137 4247 8719
 47
 846 1093 1434 1851 3258 2396 4330 4133 3629
 56
 606 969 1873 1646 1892 1861 2311 3914 3609
 33
 367 550 1241 1323 1098 1943 2662 4431 3920
 19
 313 724 703 1301 1477 1683 2457 3761 1875
 34
 187 370 529 734 970 1106 2020 2597 1879
 16
 178 296 378 754 662 1040 1335 1664 2046
 12
 5
 138 235 427 290 444 588 791 865 855
 1
 1
 27
 78
 N/K
 5
 N/K
 100 2
 80 80 80 85 80 84 100100100 87
 Run time (s) of EM algorithm
 Run time (s) of spectral algorithm
 100 40 100 80 100100100100 80 100100
 N/K
 N/K
 Accuracy (%) of spectral algorithm 100 40 82 100100100100100 75 100100
 56
 34
 10 31
 54
 89
 73
 165 172 229 266 233 216
 116 T
 158
 200

Figure 4: Clustering accuracy and run time results for synthetic data experiments. Table 1: Clustering accuracies for handwritten digit dataset. Algorithm

1v2
 1v3

1v4		
2v3		
2v4		
2v5		
Spectral EM init.	w/ Spectral EM init.	at Random
100	100	96
70	99	99
54	100	98
83	96	83
99	100	100
99	100	100

regions where we have less data, spectral algorithm is much faster. Note that, in spectral algorithm we include the time spent in moment computation. We used four restarts for EM, and take the result with highest likelihood, and used an automatic stopping criterion. 4.3

Real data experiment

We ran an experiment on the handwritten character trajectory dataset from the UCI machine learning repository [8]. We formed pairs of characters and compared the clustering results for three algorithms: the proposed spectral learning approach, EM initialized at random, and EM initialized with MoM algorithm as explored in [9]. We take the maximum accuracy of EM over 5 random initializations in the third row. We set the algorithm parameters to $K = 2$ and $M = 4$. There are 140 sequences of average length 100 per class. In the original data, $L = 3$, but to apply MoM learning, we require that $M \leq L$. To achieve this, we transformed the data vectors with a cubic polynomial feature transformation such that $L = 10$ (this is the same transformation that corresponds to a polynomial kernel). The results from these trials are shown in Table 1. We can see that although spectral learning doesn't always surpass randomly initialized EM on its own, it does serve as a very good initialization scheme.

5

Conclusions and future work

We have developed a method of moments based algorithm for learning mixture of HMMs. Our experimental results show that our approach is computationally much cheaper than EM, while being comparable in accuracy. Our real data experiment also show that our approach can be used as a good initialization scheme for EM. As future work, it would be interesting to apply the proposed approach on other hierarchical latent variable models. Acknowledgements: We would like to thank Taylan Cemgil, David Forsyth and John Hershey for valuable discussions. This material is based upon work supported by the National Science Foundation under Grant No. 1319708.

2 References

- [1] A. Anandkumar, D. Hsu, and S.M. Kakade. A method of moments for mixture models and hidden markov models. In COLT, 2012. [2] A. Anandkumar, R. Ge, D. Hsu, S.M. Kakade, and M. Telgarsky. Tensor decompositions for learning latent variable models. arXiv:1210.7559v2, 2012. 8
- [3] Daniel Hsu, Sham M. Kakade, and Tong Zhang. A spectral algorithm for learning hidden markov models a spectral algorithm for learning hidden markov models. Journal of Computer and System Sciences, (1460-1480), 2009. [4] P. Smyth. Clustering sequences with hidden markov models. In Advances in neural information processing systems, 1997. [5] Yuting Qi, J.W. Paisley, and L. Carin. Music analysis using hidden markov mixture models. Signal Processing, IEEE Transactions on, 55(11):5209–5224, nov. 2007. [6] A. Jonathan, S. Sclaroff, G. Kollios, and V. Pavlovic. Discovering clusters in motion time-series data. In CVPR, 2003. [7] Tim Oates, Laura Firoiu, and Paul R. Cohen. Clustering time series with hidden markov models and dynamic time warping. In In Proceedings of the IJCAI-99 Workshop on Neural, Symbolic and Reinforcement Learning Methods for Sequence Learning, pages 17–21, 1999. [8] K. Bache and M. Lichman. UCI machine learning repository, 2013. [9] Arun Chaganty and Percy Liang. Spectral experts for estimating mixtures of linear regressions. In International Conference on Machine Learning (ICML), 2013.