# A Spectral Algorithm for Learning Hidden Markov Models

By : Mohsen Nabian

## Summary :

This paper provides efficient algorithms for learning hidden markov models (HMMs)$_\lambda$. HMMs are generally Computationally hard to be learnt from Data and sometimes require to use heuristics which themselves makes optimization to trap in local optimum.

However this paper with the algorithm proposed, with the constraint on the problems to be in "natural separation condition" learning procedure is simple. it employs only a Singular value decomposition and matrix multiplications.

Also the complexity of the algorithm implicitly depends number of observation through spectral properties of underlying HMM.

1/

# Hidden Markov Model

HMM defined as a sequences of hidden states $(h_t)$ and observations $(x_t)$.

~~hidden states: $h_1, h_2 \ldots h_t$~~

Set of hidden states : $[m] = \{1, 2, \ldots, m\}$ $\qquad n \geqslant m$

" " observations : $[n] = \{1, 2, \ldots, n\}$

$$T_{ij} = Pr\left[h_{t+1} = i \mid h_t = j\right] \qquad T \in \mathbb{R}^{m \times m}$$

$$O_{ij} = Pr\left[x_t = i \mid h_t = j\right] \qquad O \in \mathbb{R}^{n \times n}$$

HMM satisfies conditional independence properties.

$\vec{\pi}$ is the initial state distribution $\vec{\pi} \in \mathbb{R}^m$

with $\overset{o}{\pi}_i = Pr\left[h_1 = i\right]$

Lemma 1) : For $x = 1, 2, \ldots, n$ define

$$A_x = T diag\left(O_{x,1}, \ldots, O_{x,m}\right) \qquad \text{ex: } A_1 = T_{diag}\left(O_{11}, O_{12} \ldots O_{1m}\right)$$

for any $t$ :

$$Pr\left[x_1, \ldots, x_t\right] = \underbrace{\vec{1}_m^T}_{\text{all ones}} A_{x_t} \cdots A_{x_1} \vec{\pi}$$

all ones
vector $\in \mathbb{R}^m$

**Assumption:**

Condition 1: $\vec{\Pi} > 0$ and $\vec{O}$ and $\vec{T}$ are rank $m$.

## Learning Model:

The goal is to derive cumulative distribution $\Pr[x_{1:t}]$ and conditional distribution $\Pr[x_t \mid x_{1:t-1}]$ for any sequence length $t$.

**definitions:**

$$[P_1]_i = \Pr[x_1 = i]$$
$$[P_{2,1}]_{ij} = \Pr[x_2 = i, x_1 = j]$$
$$[P_{3,x,1}]_{ij} = \Pr[x_3 = i, x_2 = x, x_1 = j] \quad \forall x \in [n]$$

where $P_1 \in \mathbb{R}^n$ is a vector, $P_{2,1} \in \mathbb{R}^{n \times n}$ and $P_{3,x,1} \in \mathbb{R}^{n \times n}$ are matrices.

Also we define $U \in \mathbb{R}^{n \times m}$ so that $U^T O$ is ⓐ invertible. an example of $U$ is 'thin' SVD of $P_{2,1}$. ↖

⇗
Condition 2

3/

$$[P_{21}]_{ij} = \sum_{k=1}^{m} \sum_{l=1}^{m} Pr[x_2 = i, x_1 = j, h_2 = k, h_1 = l]$$

$$= \sum_{k=1}^{m} \sum_{l=1}^{m} O_{ik} T_{kl} \vec{\pi}_l [\vec{O}^T]_{lj}$$

equally by Lemma 2

or in general:

$$P_{2,1} = \textcircled{O T} \, O \, T \, diag(\vec{\pi}) \, O^T$$

$$O = P_{2,1} \left( T \, diag(\vec{\pi}) \, O^T \right)^{+}$$

$\textcircled{X}^{+}$ denots the Moore–Penrose pseudo-inverse of Matrix $X$

So: $\begin{cases} range(O) \subseteq range(P_{2,1}) \\ rank(P_{2,1}) = rank(O) = m \\ range(U) = range(P_{2,1}) = range(O) \end{cases}$

So in each iteration:

→ Compute $SVD(P_{21}) \Longrightarrow$ discover $U$ that satisfies Condition 2

$\Longrightarrow \textcircled{②}$ we define $b$, $B$

$\begin{cases} \vec{b}_1 = U^T P_1 \\ \vec{b}_\infty = (P_{2,1}^T U)^+ P_1 \\ B_x = (U^T P_{3,x,1})(U^T P_{21})^+ \quad \forall x \in [n] \end{cases}$

4/

1) Spectral Learning of Hidden Markov Models :

Algorithm :

we want to predict probability of a sequence:

$$Pr\left[x_1, \ldots, x_t\right] = \hat{b}_\infty^T \hat{B}_{x_t} \ldots \hat{B}_{x_1} \hat{b}_1 .$$

given observation $x_t$, the 'internal state' update is:

$$\hat{b}_{t+1} = \frac{\hat{B}_{x_t} \hat{b}_t}{\hat{b}_\infty^T \hat{B}_{x_t} \hat{b}_t}$$

So Algorithm :    Learn HMM $(m, N)$

input : $m$ - number of states , $N$ - sample size
Return : HMM model parametrized by $\{\hat{b}_1, \hat{b}_\infty, \hat{B}_x \forall x \in [n]\}$

① Independently sample $N$ observation triples $(x_1, x_2, x_3)$ from HMM to form emperical estimates :
   $$\hat{P}_1, \hat{P}_{2,1}, \hat{P}_{3,x,1} \forall x \in [n] \text{ of } P_1, P_{2,1}, P_{3,x,1} \forall x \in [n].$$

② Compare the SVD of $\hat{P}_{2,1}$ and let $\hat{U}$ be the matrix of left singular vectors corresponding to the $m$ largest singular values

③ Compute model parameters:

(a) $\hat{b}_1 = \hat{U}^T \hat{P}_1$

(b) $\hat{b}_\infty = (\hat{P}_{2,1}^T \hat{U})^+ \hat{P}_1$

(c) $\hat{B}_x = \hat{U}^T P_{3,x,1} (\hat{U}^T \hat{P}_{2,1})^+ \quad \forall x \in [n]$

---

Now to predict Conditional Probability:

$$\hat{P}_r\left[x_t \mid x_{1:t-1}\right] = \frac{\hat{b}_\infty^T \hat{B}_{x_t} \hat{b}_t}{\sum_x \hat{b}_\infty^T \hat{B}_x \hat{b}_t}$$

$\hookrightarrow$ Conditional Probability

and remineliy

$$\hat{P}_r\left[x_1, \cdots, x_t\right] = \hat{b}_\infty^T \hat{B}_{x_t} \cdots \hat{B}_{x_1} \hat{b}_1$$

$\hookrightarrow$ joint probability

and at the end the author provides ~~the~~
Some $\underset{\text{upper}}{\wedge}$ bounds for the Learning error ($\overset{\sim}{\varepsilon}$ accuracy)

· for the case of Joint probability and Conditional Probability

6