

Computational Mechanics:

2. Hidden Markov Models

ICLAD INTENSIVE - SPRING 2026

Outline:

- 1) What is a hidden Markov Model (HMM)?
- 2) HMM minimality & uniqueness
- 3) Generalised hidden Markov Model (GHMM)

What is a hidden Markov Model?

Intuitively, a hidden Markov Model (HMM) is a stochastic process that consists of two components:

- * A 'hidden' unobserved component that evolves according to a Markovian dynamic — the next hidden state only depends on the current hidden state.
- * An observed component that depends only on the current hidden state in a possibly lossy way.

What is a hidden Markov Model?

Hidden Dynamic

Observations



What is a hidden Markov Model?

A (Mealy-type) HMM consists of:

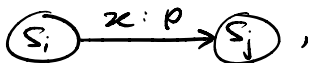
* A set S ^{"think"} ~~is~~ collection of hidden states

* A set X ^{"think"} ~~is~~ vocabulary of observations

* A collection of transition matrices $(T^{(x)})_{x \in X}$ ^{"think"} ~~is~~ dynamic that determines emissions & hidden state transitions

where $T_{ij}^{(x)} := P(X=x, s_j=s_j | s_i=s_i) = P(x, s_j | s_i)$

It is convenient to use the following graphical notation to visualise process transitions:



$$p = P(x, s_j | s_i)$$

Fix an initial distribution over hidden states

$$\eta_i^{(\phi)} := (P(s_i))_i, \quad i = 1, \dots, |S|$$

The probability of observing the sequence of emissions $w := w_1, w_2, \dots, w_n \in \mathcal{X}^n$ factorises as:

$$P(w) = \sum_{s_0, \dots, s_n} P(s_0) P(w_1, s_1 | s_0) P(w_2, s_2 | s_1) \dots P(w_n, s_n | s_{n-1})$$

Conveniently, the above expression can be stated in terms of HMM transition matrices as:

$$P(w) = \eta^{(\phi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)} \mathbf{1}$$

where $\eta^{(\phi)}$ - row vector w.r.t ordered basis of S & $\mathbf{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$

Example 1: (Biased coin)

$$S = \{s_0\}$$

$$\mathcal{X} = \{0, 1\}$$



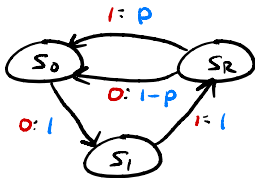
$$T^{(0)} = 1-p$$

$$T^{(1)} = p$$

Example 2: (zero-one-random)

$$S = \{s_0, s_1, s_R\}$$

$$\mathcal{X} = \{0, 1\}$$



$$T^{(0)} = \begin{array}{c} \begin{array}{ccc} s_0 & s_1 & s_R \\ \hline \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 1-p & 0 & 0 \end{bmatrix} \end{array} \begin{array}{l} s_0 \\ s_1 \\ s_R \end{array} \end{array}$$

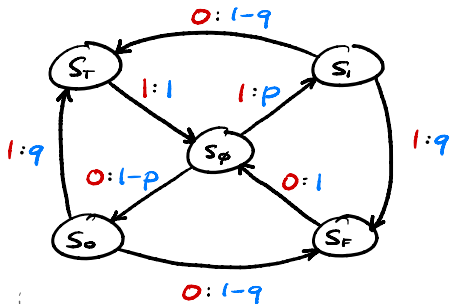
$$T^{(1)} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ p & 0 & 0 \end{bmatrix}$$

... 01R01... $R \in \{0, 1\}$

Example 3: (random-random - XOR)

$$S = \{s_\phi, s_0, s_1, s_F, s_T\}$$

$$\chi = [0, 1]$$



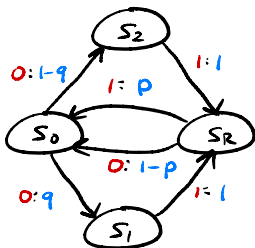
$$T^{(0)} = \begin{array}{c|ccccc} & s_\phi & s_0 & s_1 & s_F & s_T \\ \hline s_\phi & 0 & 1-p & 0 & 0 & 0 \\ \hline s_0 & 0 & 0 & 0 & 1-q & 0 \\ \hline s_1 & 0 & 0 & 0 & 0 & 1-q \\ \hline s_F & 1 & 0 & 0 & 0 & 0 \\ \hline s_T & 0 & 0 & 0 & 0 & 0 \end{array}$$

$$T^{(1)} = \begin{array}{c|ccccc} & s_\phi & s_0 & s_1 & s_F & s_T \\ \hline s_\phi & 0 & 0 & p & 0 & 0 \\ \hline s_0 & 0 & 0 & 0 & 0 & q \\ \hline s_1 & 0 & 0 & 0 & q & 0 \\ \hline s_F & 0 & 0 & 0 & 0 & 0 \\ \hline s_T & 1 & 0 & 0 & 0 & 0 \end{array}$$

Example 4:

$$S = \{s_0, s_1, s_2, s_R\}$$

$$\chi = [0, 1]$$



$$\tilde{T}^{(0)} = \begin{array}{c|cccc} & s_0 & s_1 & s_2 & s_R \\ \hline s_0 & 0 & q & 1-q & 0 \\ s_1 & 0 & 0 & 0 & 0 \\ s_2 & 0 & 0 & 0 & 0 \\ s_R & 1-p & 0 & 0 & 0 \end{array}$$

$$\tilde{T}^{(1)} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ p & 0 & 0 & 0 \end{bmatrix}$$

Example 4:

$$S = \{s_0, s_1, \dots, s_r\}$$



$$X = [0, 1, \dots]$$

Same process as
zero-one-random!

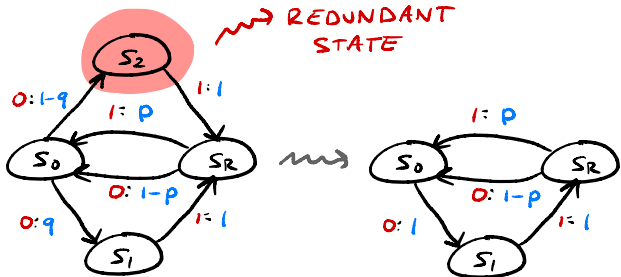
$$\tilde{T}^{(0)} = \begin{array}{c} s_0 \\ \left[\begin{array}{cccc|cc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1-p & 0 & 0 & 0 & p & 0 \end{array} \right] \end{array}$$

s_2
 s_2

Example 4:

$$S = \{s_0, s_1, s_2, s_R\}$$

$$\mathcal{X} = \{0, 1\}$$



$$\tilde{T}^{(0)} = \begin{array}{c|cccc} & s_0 & s_1 & s_2 & s_R \\ \hline s_0 & 0 & q & 1-q & 0 \\ s_1 & 0 & 0 & 0 & 0 \\ s_2 & 0 & 0 & 0 & 0 \\ s_R & 1-p & 0 & 0 & 0 \end{array}$$

$$\tilde{T}^{(1)} = \begin{array}{c|cccc} & s_0 & s_1 & s_2 & s_R \\ \hline s_0 & 0 & 0 & 0 & 0 \\ s_1 & 0 & 0 & 0 & 1 \\ s_2 & 0 & 0 & 0 & 1 \\ s_R & p & 0 & 0 & 0 \end{array}$$

Let $(T^{(x)})_{x \in \mathcal{X}} - \mathbb{Z}^1 \mathbb{R}$

$(\tilde{T}^{(x)})_{x \in \mathcal{X}} - \mathbb{Z}^1 \mathbb{R}$ with redundancy

There exist matrices U & V such that:

$$\tilde{T}^{(x)} = V T^{(x)} U \quad \& \quad Id_{3 \times 3} = UV$$

$\eta^{(\phi)} U$ - stochastic & $V T$ - all ones

for all $x \in \{0, 1\}$. It follows that

$$\eta^{(\phi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)} T$$

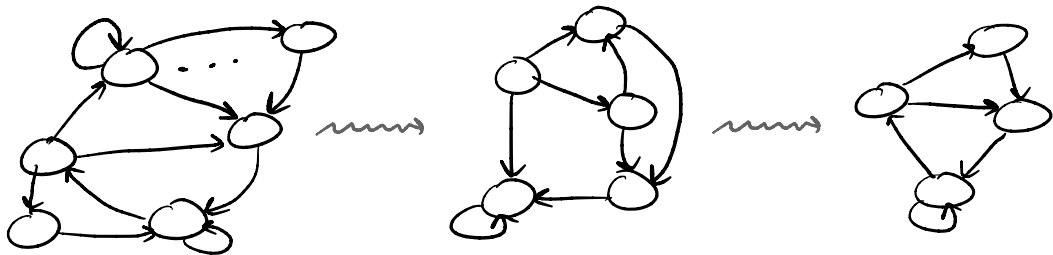
$$= \eta^{(\phi)} U V T^{(w_1)} U V T^{(w_2)} U \dots V T^{(w_n)} U V T$$

$$= \tilde{\eta}^{(\phi)} \tilde{T}^{(w_1)} \tilde{T}^{(w_2)} \dots \tilde{T}^{(w_n)} \tilde{T}$$

HMM minimality & uniqueness

A given HMM realisation of a stochastic process may be needlessly large i.e. have more hidden states than necessary.

We say that a HMM realisation of a stochastic process is minimal if there exists no other HMM realisations with fewer hidden states.



Claim: For a stochastic process that admits an HMM realisation, the minimal HMM is not unique.

(i.e. there does not exist similarity transforms relating the minimal HMMs.)

This is quite dissatisfying!

Suppose we trained a transformer on a process whose minimal HMM is not unique; what realisation would it learn the details of?

Can we identify a simple generalisation of HMMs such that uniqueness is guaranteed for the more general class?

Consider the output distribution of an HMM with $(T^{(x)})_{x \in \mathcal{X}}$ & initial hidden state dist. $\eta^{(\phi)}$

$$\begin{aligned}
 P(w) &= \eta^{(\phi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)} \mathbf{1} \\
 &= \eta^{(\phi)} \underbrace{V^{-1} V}_{I_d} T^{(w_1)} \underbrace{V^{-1} V}_{I_d} T^{(w_2)} \dots \underbrace{V^{-1} V}_{I_d} T^{(w_n)} \underbrace{V^{-1} V}_{I_d} \mathbf{1} \\
 &= \tilde{\eta}^{(\phi)} \tilde{T}^{(w_1)} \tilde{T}^{(w_2)} \dots \tilde{T}^{(w_n)} \tilde{\mathbf{1}}, \quad V \in GL(S)
 \end{aligned}$$

where $\tilde{T}^{(x)} = V T^{(x)} V^{-1}$ & $\tilde{\eta}^{(\phi)} = \eta^{(\phi)} V^{-1}$ & $\tilde{\mathbf{1}} = V \mathbf{1}$

An identical process is produced by the transition matrices $(\tilde{T}^{(x)})_{x \in \mathcal{X}}$ and the vectors $\tilde{\eta}^{(\phi)}$ & $\tilde{\mathbf{1}}$.

\Rightarrow Why constrain the transition matrix to have probabilities as matrix elements?

Generalised hidden Markov Model

A generalised hidden Markov Model (GHMM) consists of:

* A set \mathcal{X}

* A collection of transition matrices $(T^{(x)})_{x \in \mathcal{X}}$ \rightsquigarrow not necessarily sub-stochastic!

* An initial row vector $\eta^{(\phi)}$ \rightsquigarrow not necessarily stochastic!

such that:

1. There exists a (column) vector φ satisfying:

a) $T\varphi = \varphi$ where $T := \sum_{x \in \mathcal{X}} T^{(x)}$

b) $\eta^{(\phi)}\varphi = 1$

2. For all $w \in \mathcal{X}^*$, $\eta^{(\phi)} T^{(w_1)} \dots T^{(w_n)} \varphi \geq 0$.

The thrust of this definition is to afford additional freedom to $(T^{(x)})_{x \in \mathcal{X}}$ and $\eta^{(\emptyset)}$ while maintaining the same expression assigning probabilities to emission sequences:

$$P(w) = \eta^{(\emptyset)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)} \varrho$$

(This fact will be explored in the exercises.)

A consequence of this freedom is that we lose connection to explicit hidden states!

Example: (Bloch Walk)

Fix $\gamma = (2\sqrt{\alpha^2 + \beta^2})^{-1}$, where $\alpha > 0$ & $\beta \in \mathbb{R}$

$$T^{(0)} = \begin{bmatrix} \frac{1}{4} & 0 & 2\alpha\beta\gamma^2 \\ 0 & (\alpha^2 - \beta^2)\gamma^2 & 0 \\ 2\alpha\beta\gamma^2 & 0 & \frac{1}{4} \end{bmatrix} \quad T^{(1)} = \begin{bmatrix} \frac{1}{4} & 0 & -2\alpha\beta\gamma^2 \\ 0 & (\alpha^2 - \beta^2)\gamma^2 & 0 \\ -2\alpha\beta\gamma^2 & 0 & \frac{1}{4} \end{bmatrix}$$

$$T^{(2)} = \begin{bmatrix} \frac{1}{4} & 2\alpha\beta\gamma^2 & 0 \\ 2\alpha\beta\gamma^2 & \frac{1}{4} & 0 \\ 0 & 0 & (\alpha^2 - \beta^2)\gamma^2 \end{bmatrix} \quad T^{(3)} = \begin{bmatrix} \frac{1}{4} & -2\alpha\beta\gamma^2 & 0 \\ -2\alpha\beta\gamma^2 & \frac{1}{4} & 0 \\ 0 & 0 & (\alpha^2 - \beta^2)\gamma^2 \end{bmatrix}$$

$$\eta^{(\emptyset)} = [1 \ 0 \ 0]$$

$$\varphi = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

Returning to our initial motivation...

Claim (Upper '97): For a stochastic process that admits a GHMM realisation, the minimal GHMM is unique (up to similarity transformation).

⇒ Transformers trained on data emitted from a GHMM should learn the structure corresponding to the minimal GHMM

We will discuss this 'structure' in the next block!

Questions?