

# Computational Mechanics:

## 3. Prediction

ILIAD INTENSIVE - SPRING 2026

## Outline:

- 1) Belief state
- 2) Next-token (probability) vector
- 3) Mixed-state presentation
- 4) GHMM prediction

## Belief state

How should we predict tokens emitted from an HMM?

↳ By def<sup>n</sup> the next-token distribution is a function of the current distribution over hidden states.

↳ Use this to make predictions

Bayes' rule tells us that the posterior prob. of occupying state  $s_j$  after observing  $x$  is:

$$P(s_j | x) = \frac{P(x, s_j)}{P(x)} = \frac{\sum_{s_0 \in S} P(s_0) P(x, s_j | s_0)}{P(x)}$$

hidden state      observed tokens

Similarly, posterior prob. of occupying  $s_j$  after observing the sequence  $W = w_1, w_2 \dots w_n$  is:

$$P(s_j | W) = \frac{\sum_{s_0, \dots, s^{n-1}} P(s_0) P(w_1, s_1 | s_0) P(w_2, s_2 | s_1) \dots P(w_n, s_j | s^{n-1})}{P(W)}$$

Expressing this object in terms of HMM data:

$$P(s_j | W) = \left( \frac{\eta(\phi) T^{(w_1)} T^{(w_2)} \dots T^{(w_n)}}{\underbrace{\eta(\phi) T^{(w_1)} T^{(w_2)} \dots T^{(w_n)}}_{\text{row-vector}}} \right)_j$$

where we order the hidden states  $S = \{s_0, s_1, \dots\}$ .

This motivates the def<sup>n</sup> of the belief state:

$$\eta^{(w)} := \frac{\eta^{(\phi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)}}{\eta^{(\phi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)} \mathbf{1}}, \quad w = w_1 w_2 \dots w_n \in \mathcal{X}^*$$

whose  $j^{\text{th}}$  element corresponds to the prob. that the hidden state is  $S_j$  after obs.  $w$ .

Suppose we observe a further token  $x \in \mathcal{X}$ , the predictive vector update rule is simple:

$$\eta^{(w)} \mapsto \eta^{(wx)} = \frac{\eta^{(w)} T^{(x)}}{\eta^{(w)} T^{(x)} \mathbf{1}}$$

## Next-token (probability) vector

For predicting next-token emissions, we actually care about:

$$P(x_j | w) = \frac{P(w x_j)}{P(w)} = \frac{\eta^{(\phi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)} T^{(x_j)} \tau}{\eta^{(\phi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)} \tau}$$

next-token      observed tokens

This motivates the def<sup>n</sup> of the next-token vector:

$$p^{(w)} := \left( \frac{\eta^{(\phi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)} T^{(x_j)} \tau}{\eta^{(\phi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)} \tau} \right)_{j=0,1,\dots}$$

where we order the tokens  $\mathcal{X} = \{x_0, x_1, \dots\}$ .

The  $j^{\text{th}}$ -element of  $p^{(w)}$  corresponds to the prob. that the next-token is  $x_j$  after obs.  $w$ .

The belief state appears directly in the def<sup>n</sup> of the next-token vector:

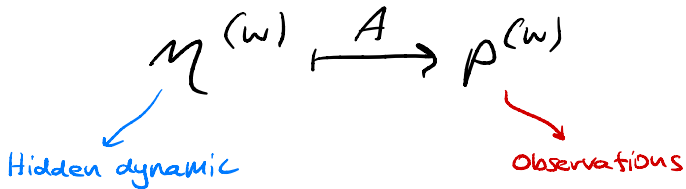
$$p_j^{(w)} := \frac{\eta^{(\varphi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)} T^{(x_j)} \gamma}{\eta^{(\varphi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)} \gamma} = \eta^{(w)} T^{(x_j)} \gamma$$

It follows that there exists a matrix that maps  $\eta^{(w)} \mapsto p^{(w)}$  i.e.:

$$p^{(w)} = \eta^{(w)} A$$

where the matrix elements of  $A$  are:

$$A_{ij} := \sum_{k=0}^{|\mathcal{S}|-1} T_{ik}^{(x_j)}, \quad \begin{array}{l} i = 0, 1, \dots, |\mathcal{S}|-1 \\ j = 0, 1, \dots, |\mathcal{X}|-1 \end{array}$$



If  $\ker(A) = \emptyset$ :



If  $\ker(A) \neq \emptyset$ :



Belief states are also useful for all future preds.

Let  $w^{(p)} = w_1^{(p)} w_2^{(p)} \dots w_m^{(p)}$  &  $w^{(f)} = w_1^{(f)} w_2^{(f)} \dots w_n^{(f)}$ :

$$\begin{aligned} P(w^{(f)} | w^{(p)}) &= \frac{\eta^{(\phi)} T^{(w_1^{(p)})} \dots T^{(w_m^{(p)})} T^{(w_1^{(f)})} \dots T^{(w_n^{(f)})} \uparrow}{\eta^{(\phi)} T^{(w_1^{(p)})} \dots T^{(w_m^{(p)})} \uparrow} \\ &= \eta^{(w^{(p)})} T^{(w_1^{(f)})} \dots T^{(w_n^{(f)})} \uparrow \end{aligned}$$

As in the next-token case there exists a multi-linear map from  $\eta^{(w^{(p)})}$  to all future predictions.

## Mixed-state presentation

For any HMM, the deterministic update rule  $\eta^{(w)} \mapsto \eta^{(wx)}$  defines a meta-dynamic over beliefs called the mixed-state presentation (MSP).

The MSP is itself a (unifilar) HMM:

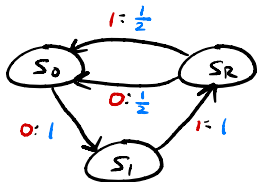


whose hidden states are the belief states.

[ For unifilar HMMs, an  $\varepsilon$ -machine is the MSP of the corresponding minimal HMM. ]

# Example 1: (zero-one-random $p = \frac{1}{2}$ )

Recall:



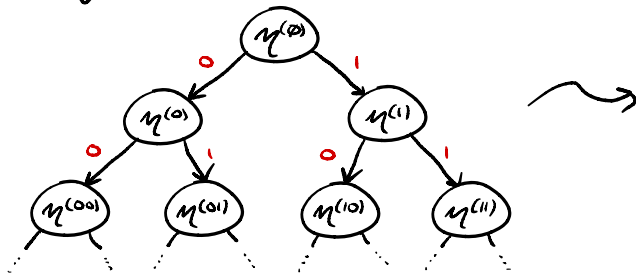
$$T^{(0)} = \begin{matrix} & \begin{matrix} S_0 & S_1 & S_R \end{matrix} \\ \begin{matrix} S_0 \\ S_1 \\ S_R \end{matrix} & \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ \frac{1}{2} & 0 & 0 \end{bmatrix} \end{matrix}$$

$$T^{(1)} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ \frac{1}{2} & 0 & 0 \end{bmatrix}$$

$$\eta^{(w)} = \frac{\eta^{(\emptyset)} T^{(w)}}{\eta^{(\emptyset)} T^{(w)} \mathbf{1}}$$

$$T^{(w)} := T^{(w_1)} \dots T^{(w_n)}$$

Deriving the MSP compute:

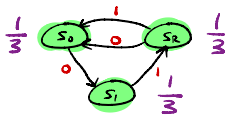


Identify nodes

where

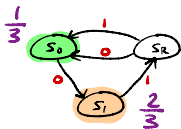
$$\eta^{(w)} = \eta^{(w')}$$

Fix:  $M^{(\emptyset)} = \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$

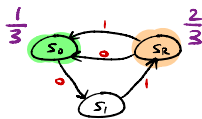


and compute:

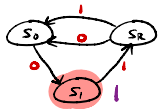
$M^{(0)} = \begin{bmatrix} \frac{1}{3} & \frac{2}{3} & 0 \end{bmatrix}$



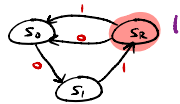
$M^{(1)} = \begin{bmatrix} \frac{1}{3} & 0 & \frac{2}{3} \end{bmatrix}$



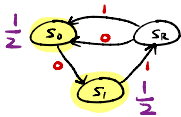
$M^{(\infty)} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$



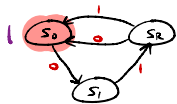
$M^{(01)} = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}$



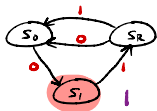
$M^{(10)} = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix}$



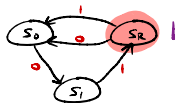
$M^{(11)} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$



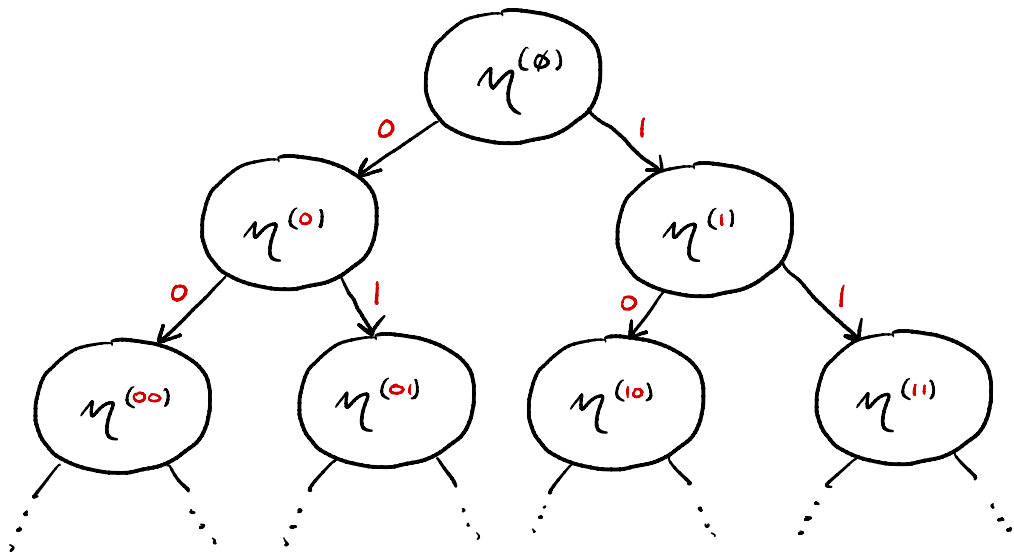
$M^{(100)} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$

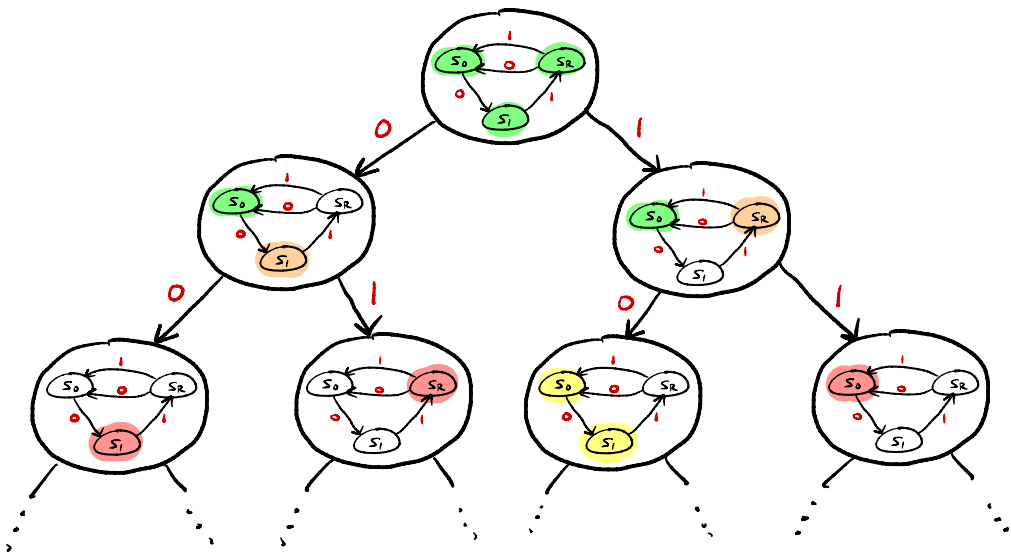


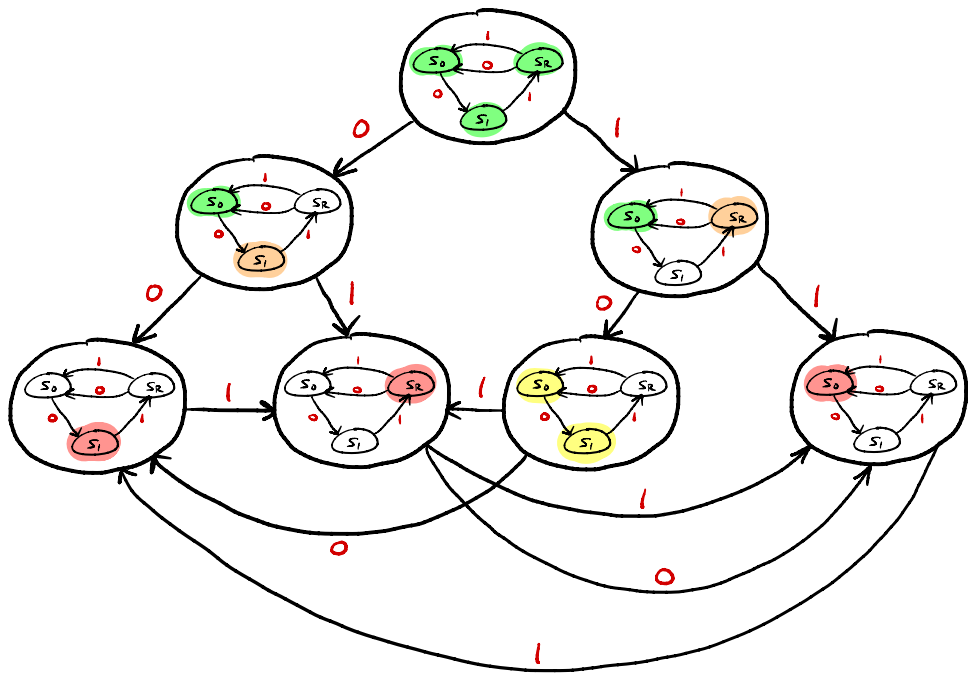
$M^{(101)} = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}$

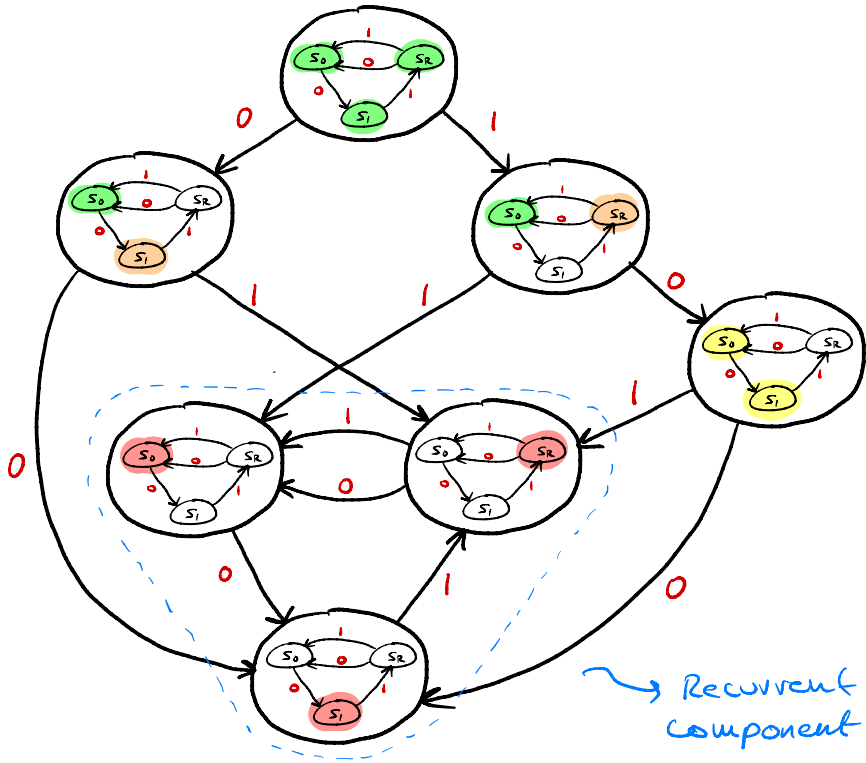


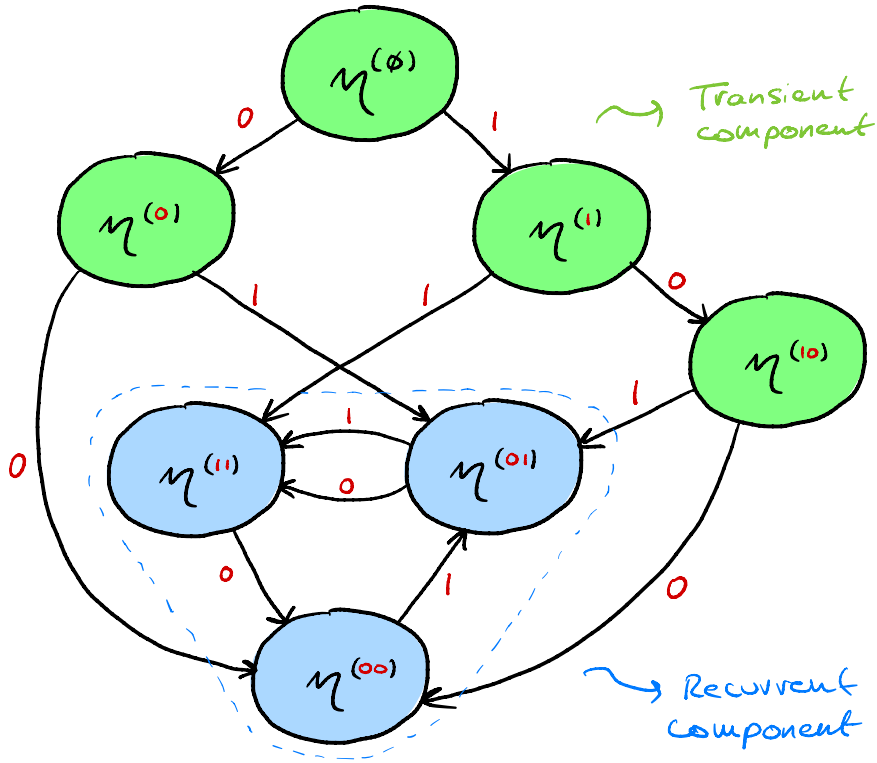
Inference process converges!







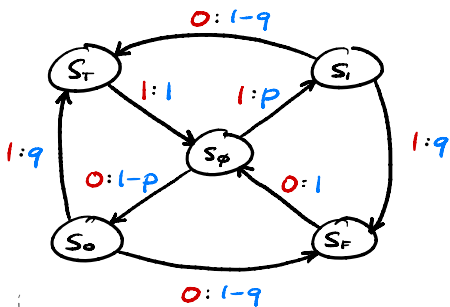




Example 2: (random-random - XOR)

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

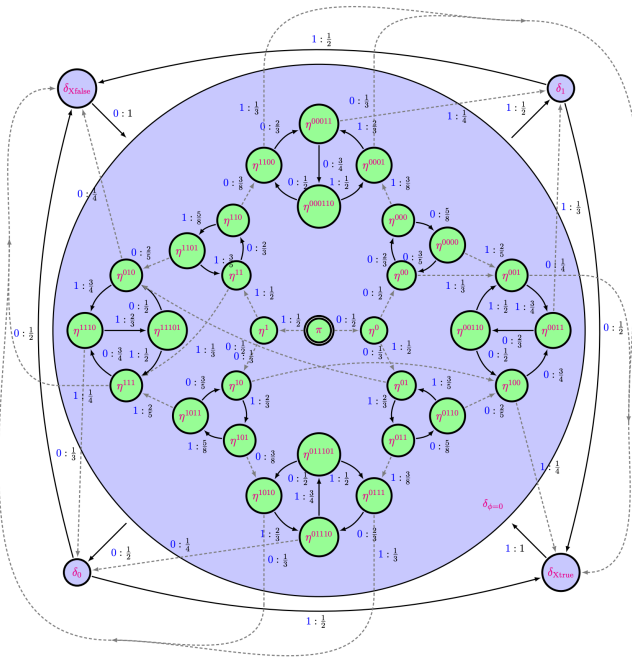
$$\chi = \{0, 1\}$$



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

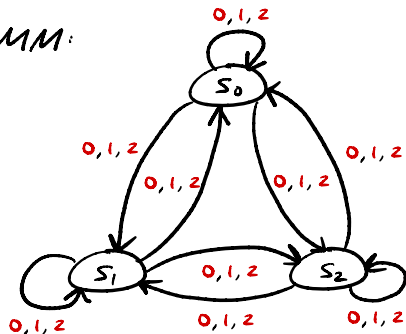
$$T^{(1)} = \begin{array}{c|ccccc} & s_\emptyset & s_0 & s_1 & s_F & s_T \\ \hline s_\emptyset & 0 & 0 & p & 0 & 0 \\ \hline s_0 & 0 & 0 & 0 & 0 & q \\ \hline s_1 & 0 & 0 & 0 & q & 0 \\ \hline s_F & 0 & 0 & 0 & 0 & 0 \\ \hline s_T & 1 & 0 & 0 & 0 & 0 \end{array}$$

The corresponding MSP:



# Example 3: (Mess 3)

HMM:

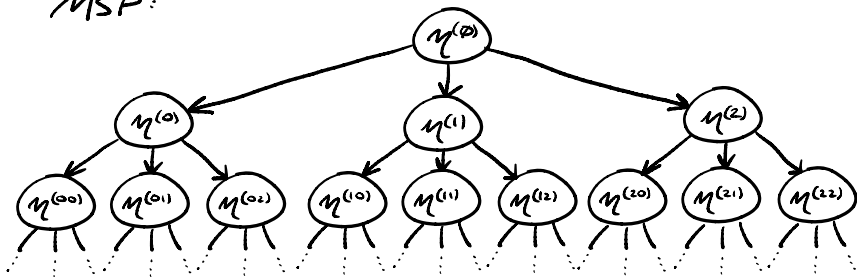


$$T^{(0)} = \begin{matrix} & \begin{matrix} S_0 & S_1 & S_2 \end{matrix} \\ \begin{matrix} S_0 \\ S_1 \\ S_2 \end{matrix} & \begin{bmatrix} \alpha y & \beta x & \beta x \\ \alpha x & \beta y & \beta x \\ \alpha x & \beta x & \beta y \end{bmatrix} \end{matrix}$$

$$T^{(1)} = \begin{bmatrix} \beta y & \alpha x & \beta x \\ \beta x & \alpha y & \beta x \\ \beta x & \alpha x & \beta y \end{bmatrix}$$

$$T^{(2)} = \begin{bmatrix} \beta y & \beta x & \alpha x \\ \beta x & \beta y & \alpha x \\ \beta x & \beta x & \alpha y \end{bmatrix}$$

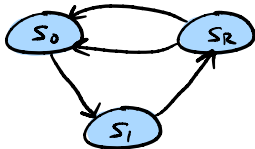
MSP:



Inference does not converge!

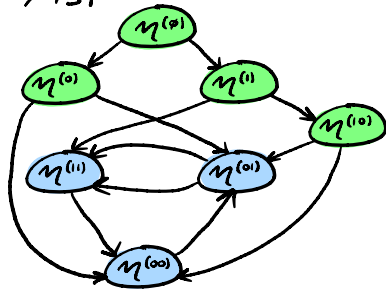
# Summary:

HMM:



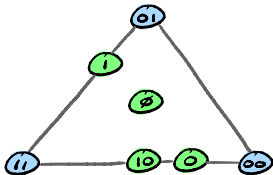
inference

MSP:



geometry

Belief-state geometry:



linear map

Next-token geometry:



## GHMM prediction:

Throughout we have considered prediction of HMM data. Fortunately, all of the relevant objects have straightforward analogues in the generalised setting.

Recall that a GHMM consists of:

$(T^{(x)})_{x \in \mathcal{X}}$       $\eta^{(\phi)}$  - initial state      $\mathcal{C}$  - final state

Belief State (aka predictive vector):

$$\eta^{(w)} := \frac{\eta^{(\phi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)}}{\eta^{(\phi)} T^{(w_1)} T^{(w_2)} \dots T^{(w_n)} \mathcal{C}}, \quad w = w_1 w_2 \dots w_n \in \mathcal{X}^*$$

Next-token vector:

$$\rho^{(w)} = (\eta^{(w)T} \langle x_j | \varphi \rangle)_{j=0,1,\dots}$$

Linear map  $\eta^{(w)} \xrightarrow{A} \rho^{(w)}$ :

$$A_{ij} := \sum_{k=0}^{d-1} T_{ik} \langle x_j | \varphi_k \rangle \quad \begin{array}{l} i = 0, 1, \dots, d-1 \\ j = 0, 1, \dots, |X|-1 \end{array}$$

A GHMM does not have 'underlying' hidden states so  $\eta^{(w)}$  is not stochastic and thus cannot be viewed as a posterior distribution.

↳ Despite this the belief state remains similarly instrumental at making token predictions.

Questions?