

# Machine Learning for Data Science (CS4786)

## Lecture 20

### Hidden Markov Models

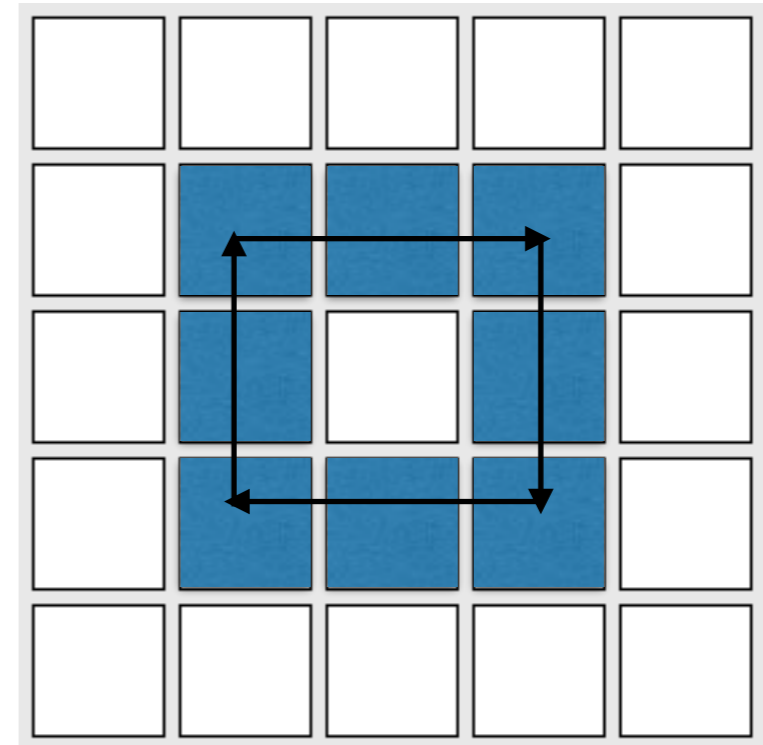
# HIDDEN MARKOV MODEL (HMM)

Same example:

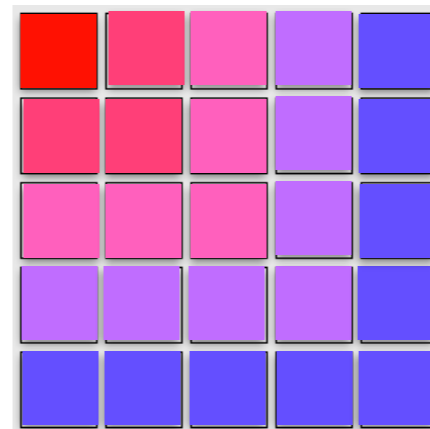


But you don't observe location  
(dark room)

You hear how close the bot is!



What you hear:

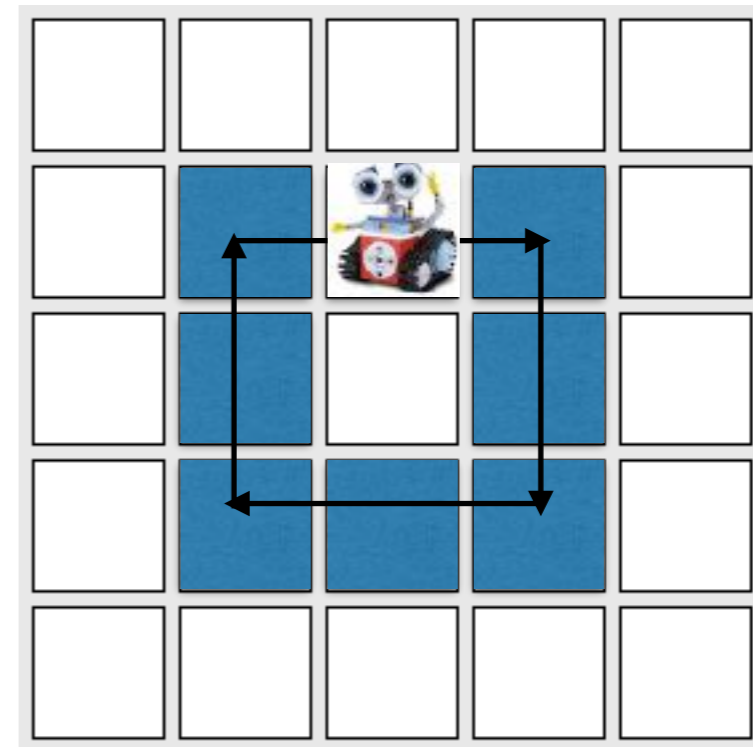


# HIDDEN MARKOV MODEL (HMM)

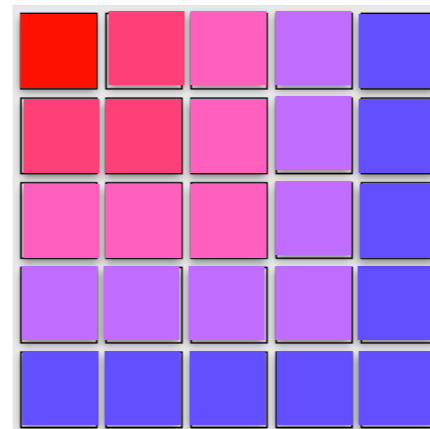
Same example:

But you don't observe location  
(dark room)

You hear how close the bot is!

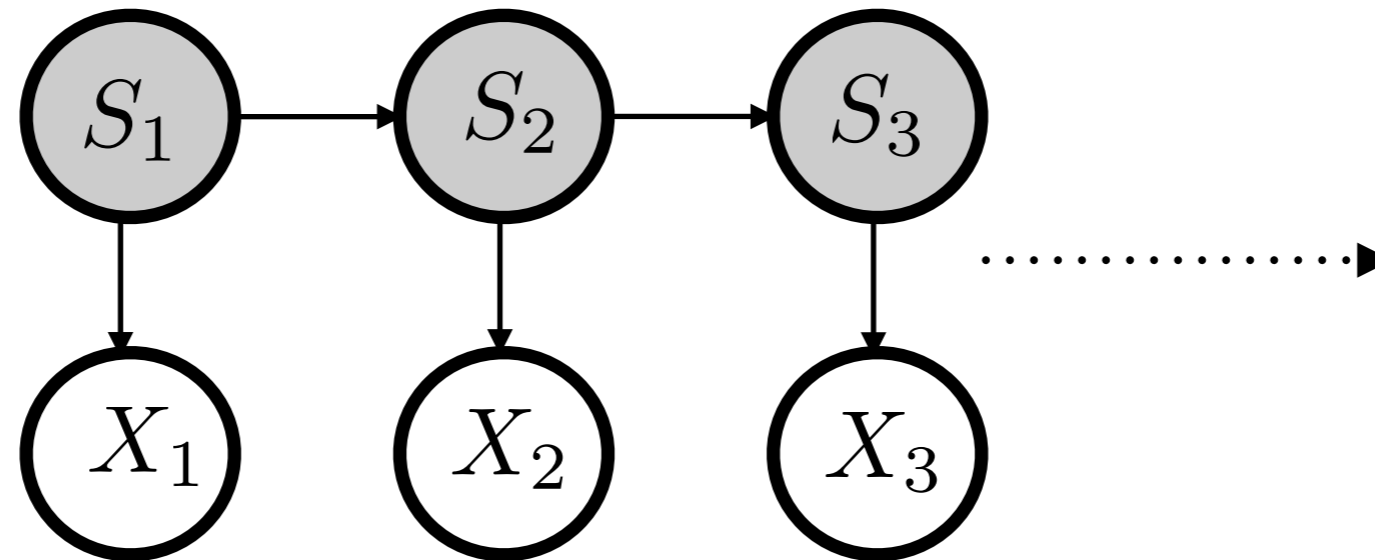


What you hear:



Can you catch the Bot?

# HIDDEN MARKOV MODEL (HMM)



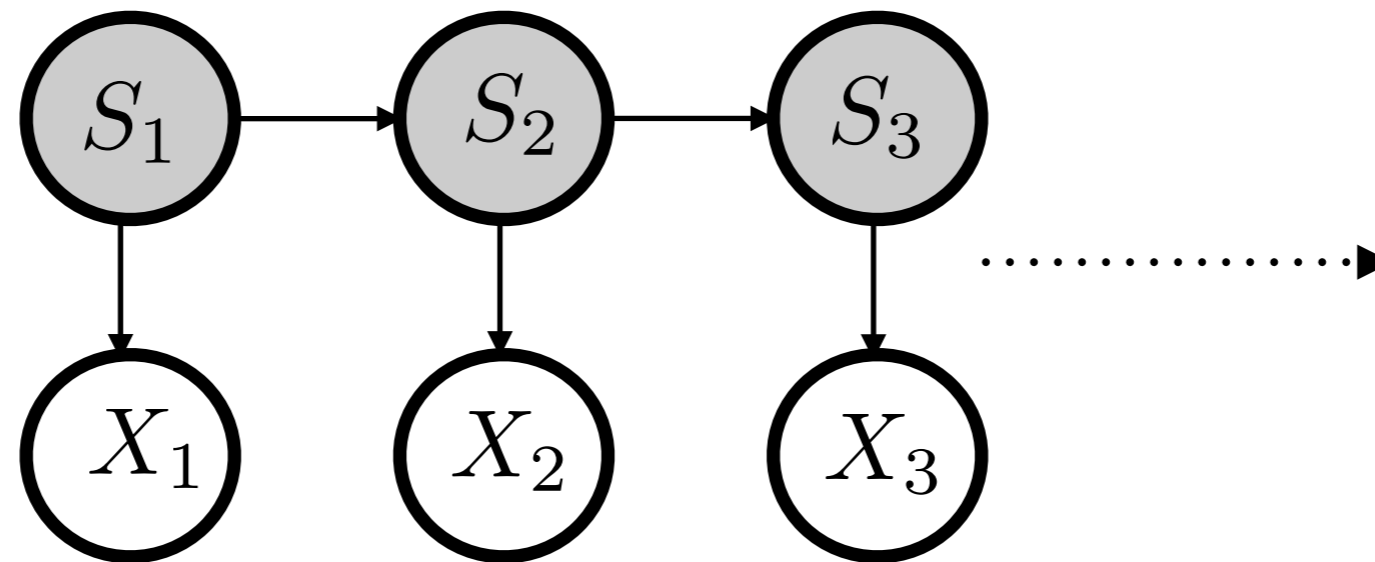
$X_t$ 's are what you hear (observation)

$S_t$ 's are the unseen locations (states)

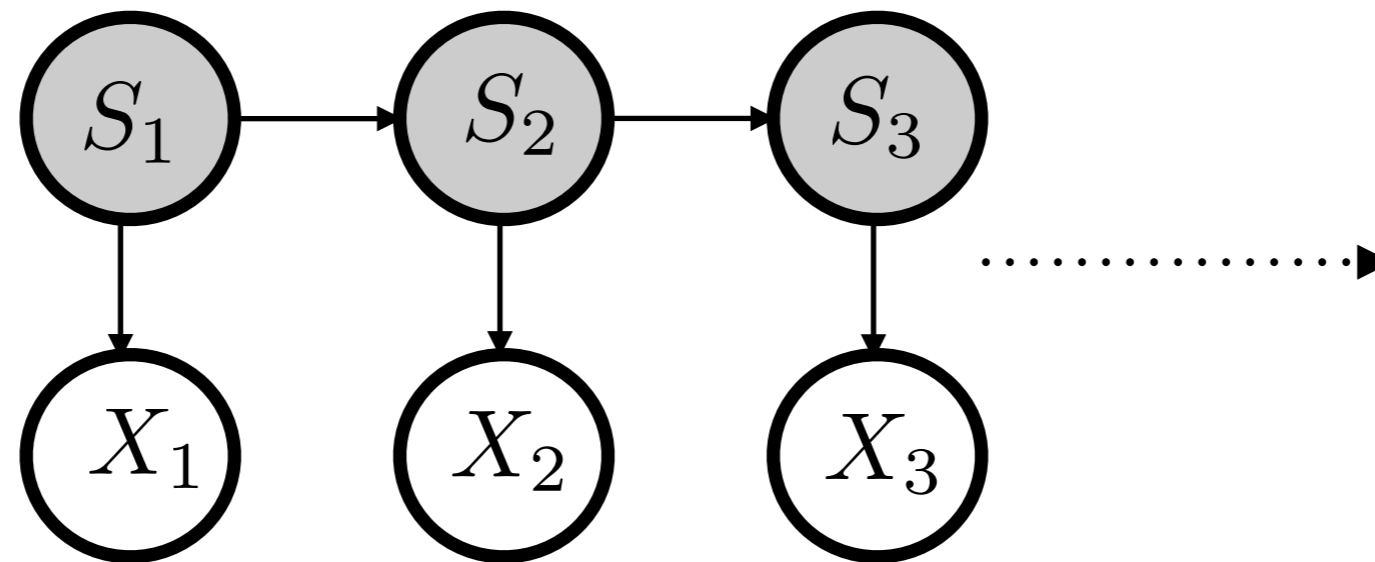
Eg: for  $n \times n$  grid we have,  $K = n^2$  states

Number of alphabets = 5  
(colors you can observe)

# HIDDEN MARKOV MODEL (HMM)

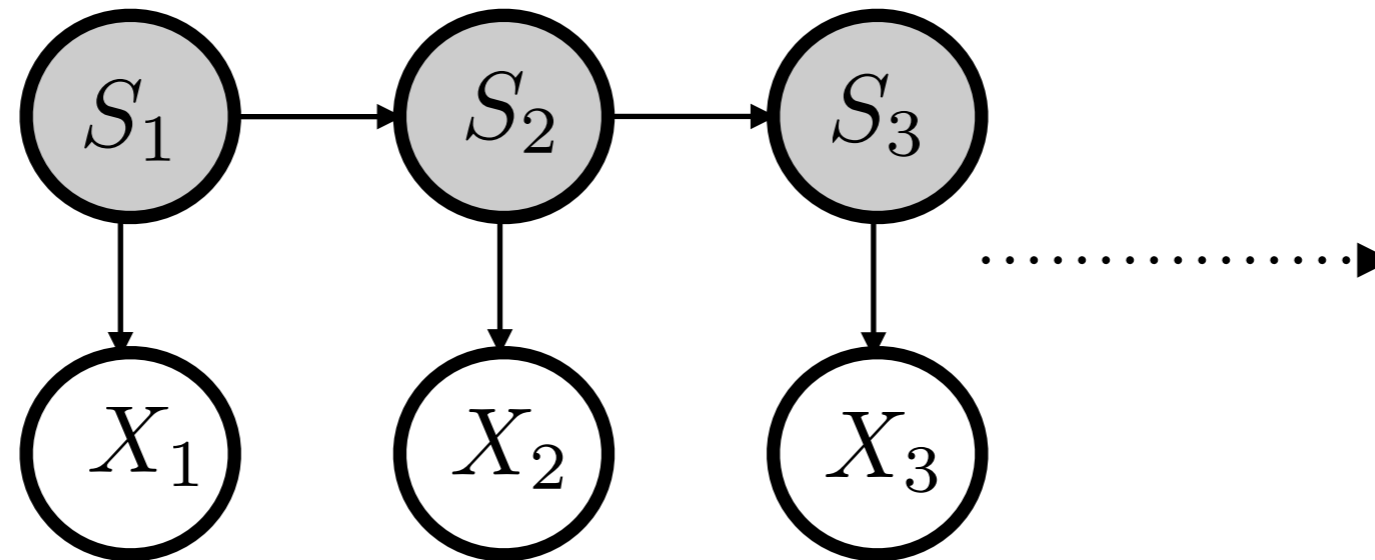


# HIDDEN MARKOV MODEL (HMM)



What are the parameters?

# HIDDEN MARKOV MODEL (HMM)



What are the parameters?

**Transition Probability table:  $T = P(S_t|S_{t-1})$**

**Emission Probabilities:  $E = P(X_t|S_t)$**

**Initial State Probabilities:  $P(S_1)$**

# HIDDEN MARKOV MODEL (HMM)



# HIDDEN MARKOV MODEL (HMM)

- What is probability that bot will be in location  $k$  at time  $t$  given the entire sequence of observations?

# HIDDEN MARKOV MODEL (HMM)

- What is probability that bot will be in location  $k$  at time  $t$  given the entire sequence of observations?

$$P(S_t = k | X_1, \dots, X_N)?$$

# INFERENCE IN HMM

$$P(S_t = k | X_1, \dots, X_N)$$

# INFERENCE IN HMM

$$P(S_t = k | X_1, \dots, X_N)$$

$$\propto P(X_{t+1}, \dots, X_N | S_t = k, X_1, \dots, X_t) P(S_t = k | X_1, \dots, X_t)$$

# INFERENCE IN HMM

$$P(S_t = k | X_1, \dots, X_N)$$

$$\propto P(X_{t+1}, \dots, X_N | S_t = k, X_1, \dots, X_t) P(S_t = k | X_1, \dots, X_t)$$

$$\propto P(X_{t+1}, \dots, X_N | S_t = k, X_1, \dots, X_t) P(S_t = k, X_1, \dots, X_t)$$

# INFERENCE IN HMM

$$P(S_t = k | X_1, \dots, X_N)$$

$$\propto P(X_{t+1}, \dots, X_N | S_t = k, X_1, \dots, X_t) P(S_t = k | X_1, \dots, X_t)$$

$$\propto P(X_{t+1}, \dots, X_N | S_t = k, X_1, \dots, X_t) P(S_t = k, X_1, \dots, X_t)$$

$$\propto P(X_{t+1}, \dots, X_N | S_t = k, X_1, \dots, X_t) P(X_t | S_t = k, X_1, \dots, X_{t-1}) P(S_t = k, X_1, \dots, X_{t-1})$$

# INFERENCE IN HMM

$$P(S_t = k | X_1, \dots, X_N)$$

$$\propto P(X_{t+1}, \dots, X_N | S_t = k, X_1, \dots, X_t) P(S_t = k | X_1, \dots, X_t)$$

$$\propto P(X_{t+1}, \dots, X_N | S_t = k, X_1, \dots, X_t) P(S_t = k, X_1, \dots, X_t)$$

$$\propto P(X_{t+1}, \dots, X_N | S_t = k, X_1, \dots, X_t) P(X_t | S_t = k, X_1, \dots, X_{t-1}) P(S_t = k, X_1, \dots, X_{t-1})$$

$$\propto P(X_{t+1}, \dots, X_N | S_t = k) P(X_t | S_t = k) P(S_t = k, X_1, \dots, X_{t-1})$$

# INFERENCE IN HMM

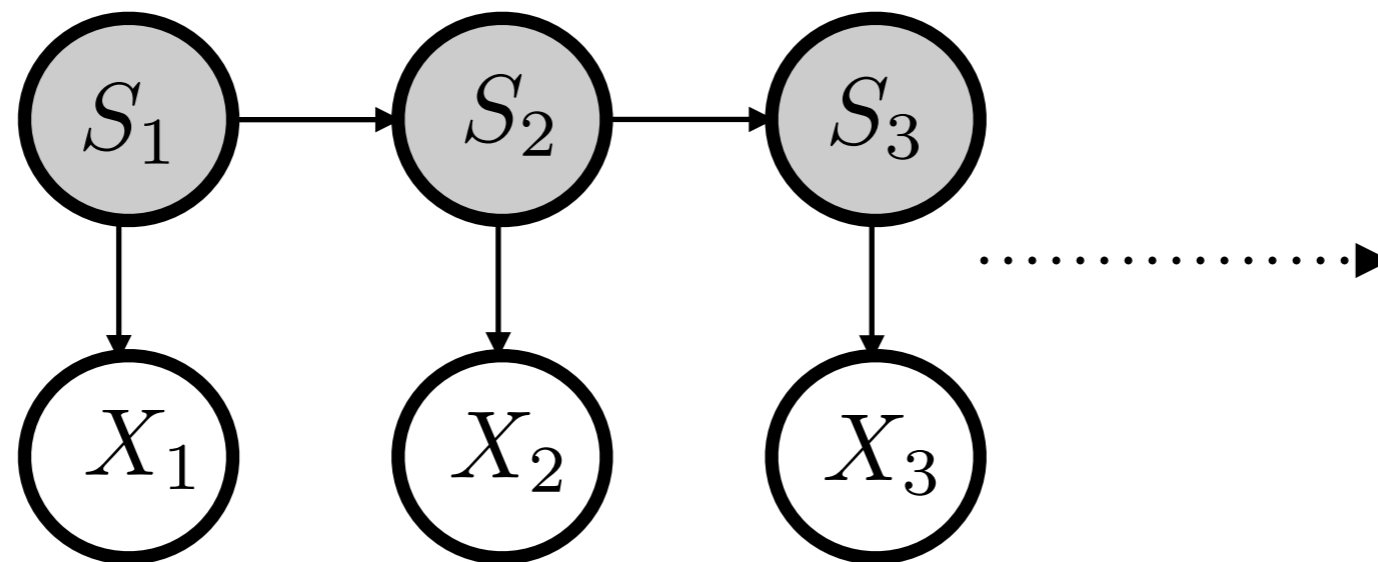
$$\begin{aligned} P(S_t = k | X_1, \dots, X_N) & \\ & \propto P(X_{t+1}, \dots, X_N | S_t = k, X_1, \dots, X_t) P(S_t = k | X_1, \dots, X_t) \\ & \propto P(X_{t+1}, \dots, X_N | S_t = k, X_1, \dots, X_t) P(S_t = k, X_1, \dots, X_t) \\ & \propto P(X_{t+1}, \dots, X_N | S_t = k, X_1, \dots, X_t) P(X_t | S_t = k, X_1, \dots, X_{t-1}) P(S_t = k, X_1, \dots, X_{t-1}) \\ & \propto P(X_{t+1}, \dots, X_N | S_t = k) P(X_t | S_t = k) P(S_t = k, X_1, \dots, X_{t-1}) \end{aligned}$$

We know  $P(X_t | S_t = k)$ 's and  $P(S_t | S_{t-1})$

Compute  $P(X_{t+1}, \dots, X_N | S_t = k)$  and  $P(S_t = k, X_1, \dots, X_t)$  recursively



# INFERENCE IN HMM

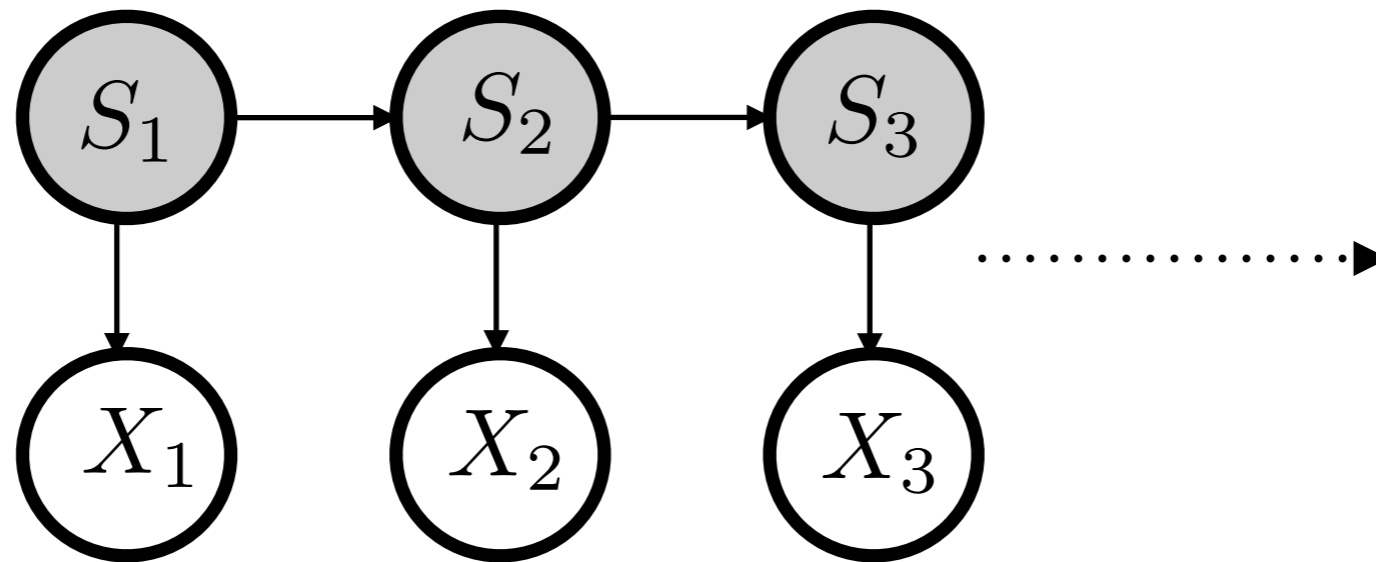


$$\text{message}_{S_{t-1} \mapsto S_t}(k) = P(S_t = k, X_1, \dots, X_{t-1})$$

$$\text{message}_{S_{t+1} \mapsto S_t}(k) = P(X_n, \dots, X_{t+1} | S_t = k)$$

$$P(S_t = k | X_1, \dots, X_n) \propto \text{message}_{S_{t-1} \mapsto S_t}(k) \times \text{message}_{S_{t+1} \mapsto S_t}(k) \times P(X_t | S_t = k)$$

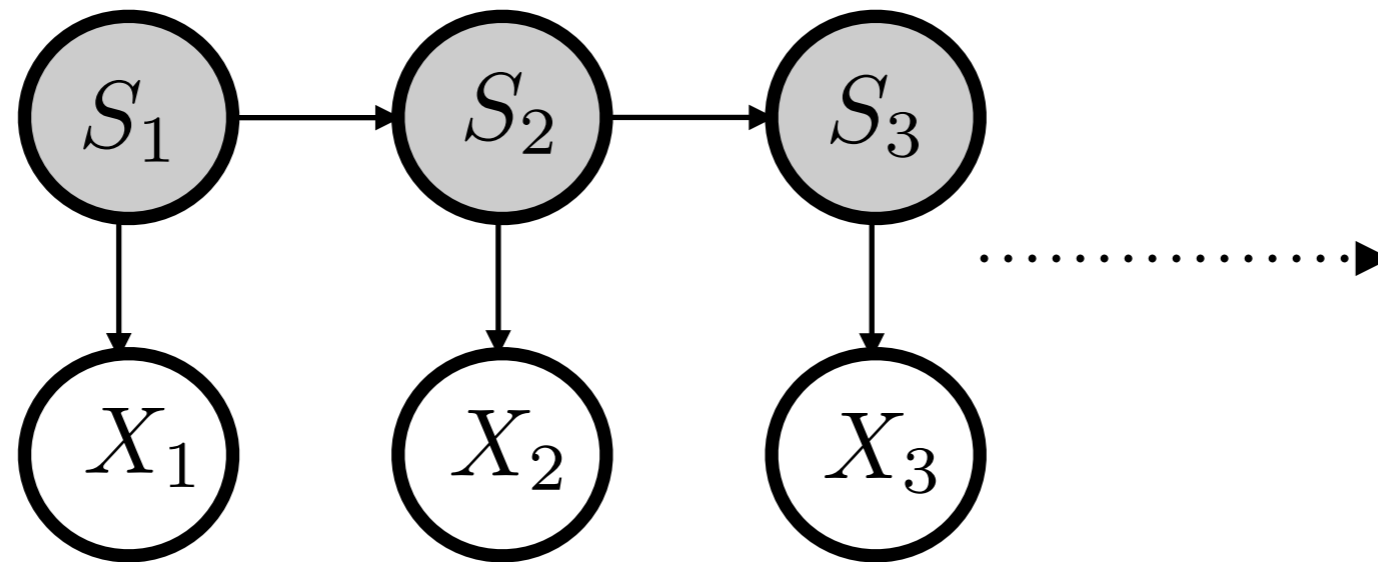
# INFERENCE IN HMM



$$\text{message}_{S_{t-1} \mapsto S_t}(k) = P(S_t = k, X_1, \dots, X_{t-1})$$

$$\text{message}_{S_{t+1} \mapsto S_t}(k) = P(X_n, \dots, X_{t+1} | S_t = k)$$

# INFERENCE IN HMM



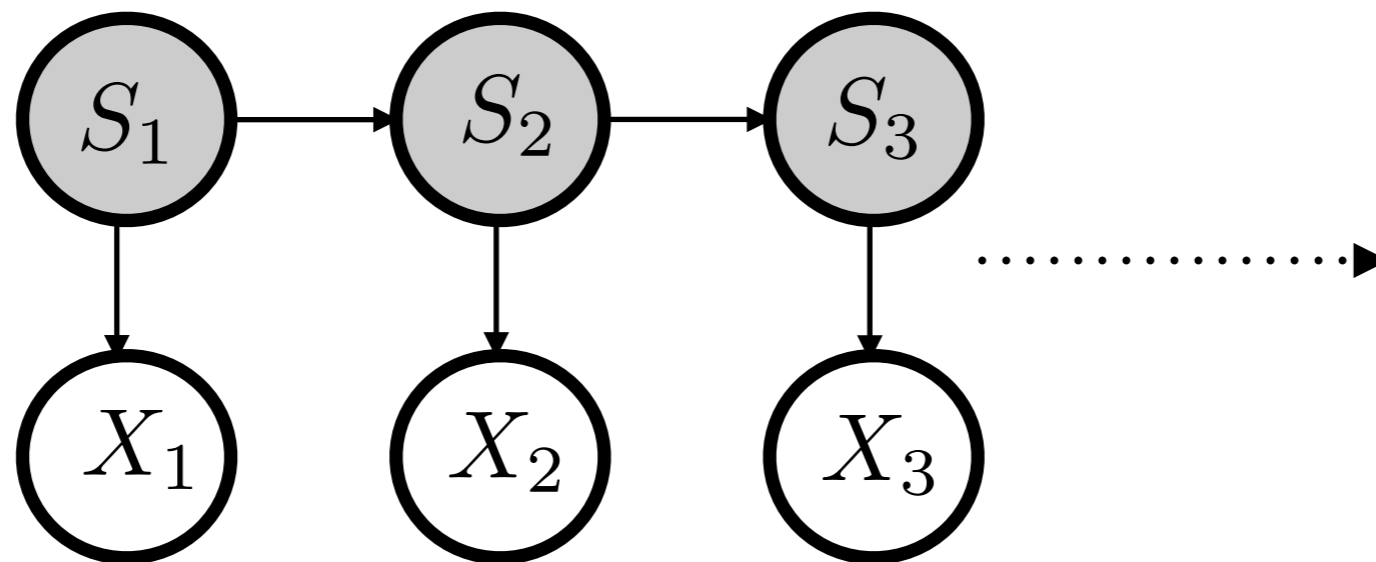
$$\text{message}_{S_{t-1} \mapsto S_t}(k) = P(S_t = k, X_1, \dots, X_{t-1})$$

$$\text{message}_{S_{t+1} \mapsto S_t}(k) = P(X_n, \dots, X_{t+1} | S_t = k)$$

Forward:

$$P(X_1, \dots, X_{t-1}, S_t = k) = \sum_{j=1}^K P(S_t = k | S_{t-1} = j) P(X_{t-1} | S_{t-1} = j) P(X_1, \dots, X_{t-2}, S_{t-1} = j)$$

# INFERENCE IN HMM



$$\text{message}_{S_{t-1} \mapsto S_t}(k) = P(S_t = k, X_1, \dots, X_{t-1})$$

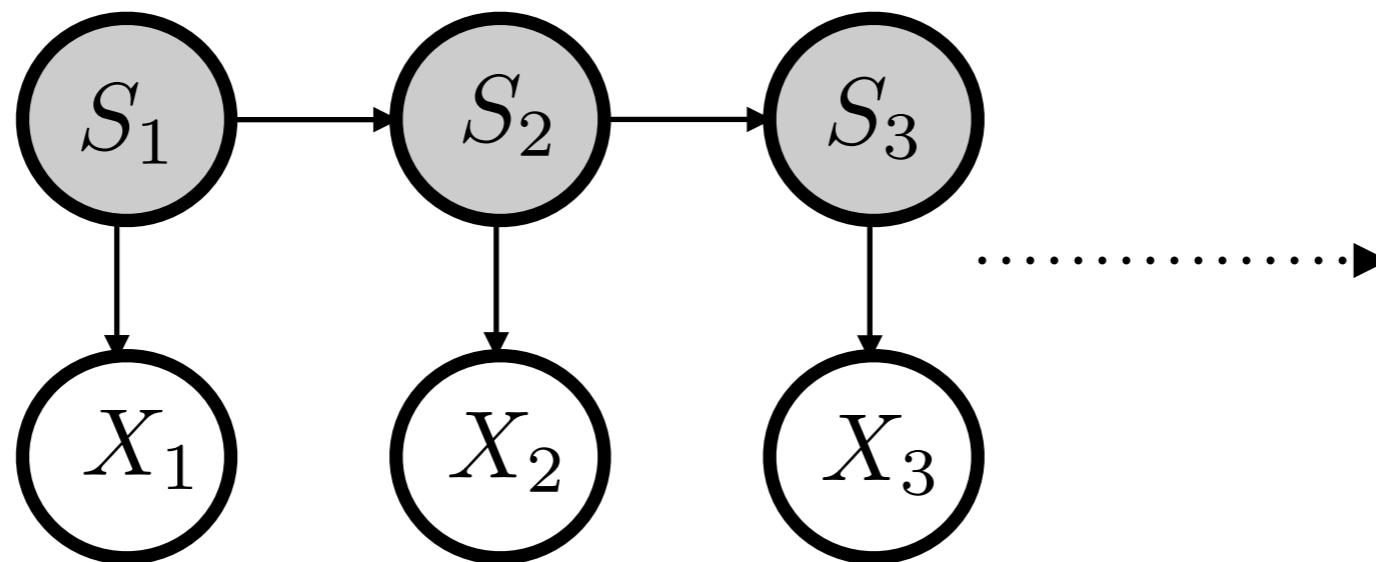
$$\text{message}_{S_{t+1} \mapsto S_t}(k) = P(X_n, \dots, X_{t+1} | S_t = k)$$

Forward:

$$P(X_1, \dots, X_{t-1}, S_t = k) = \sum_{j=1}^K P(S_t = k | S_{t-1} = j) P(X_{t-1} | S_{t-1} = j) P(X_1, \dots, X_{t-2}, S_{t-1} = j)$$

$$\text{message}_{S_{t-1} \mapsto S_t}(k) = \sum_{j=1}^K P(S_t = k | S_{t-1} = j) P(X_{t-1} | S_{t-1} = j) \text{message}_{S_{t-2} \mapsto S_{t-1}}(j)$$

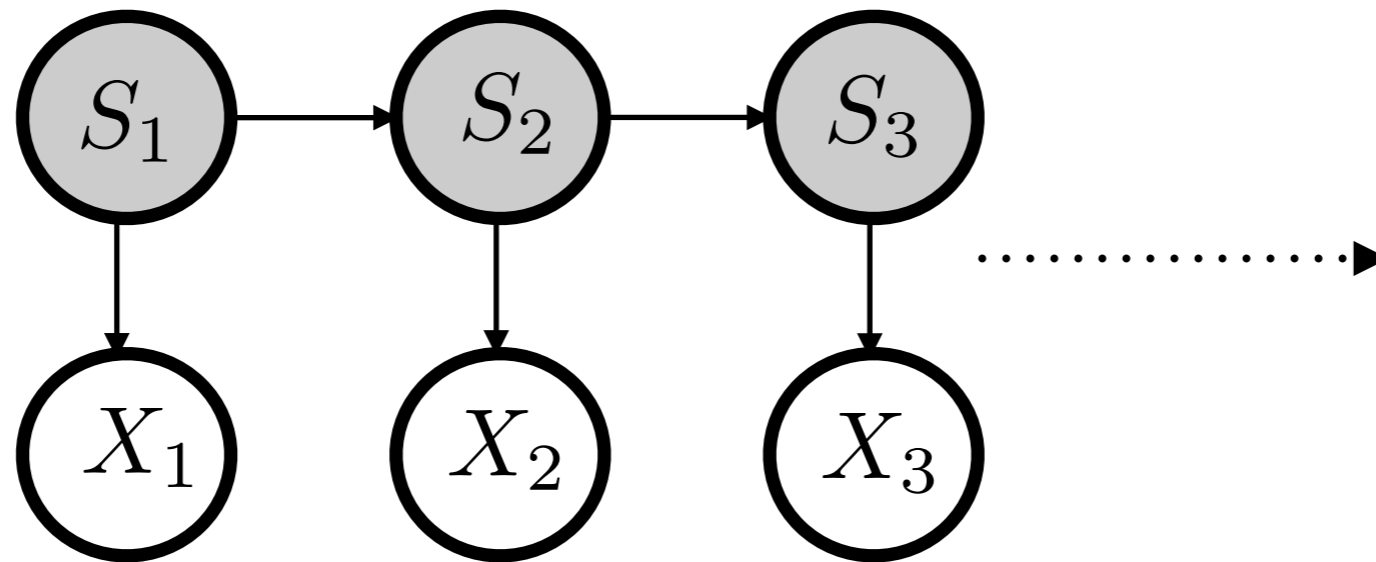
# INFERENCE IN HMM



$$\text{message}_{S_{t-1} \mapsto S_t}(k) = P(S_t = k, X_1, \dots, X_{t-1})$$

$$\text{message}_{S_{t+1} \mapsto S_t}(k) = P(X_n, \dots, X_{t+1} | S_t = k)$$

# INFERENCE IN HMM



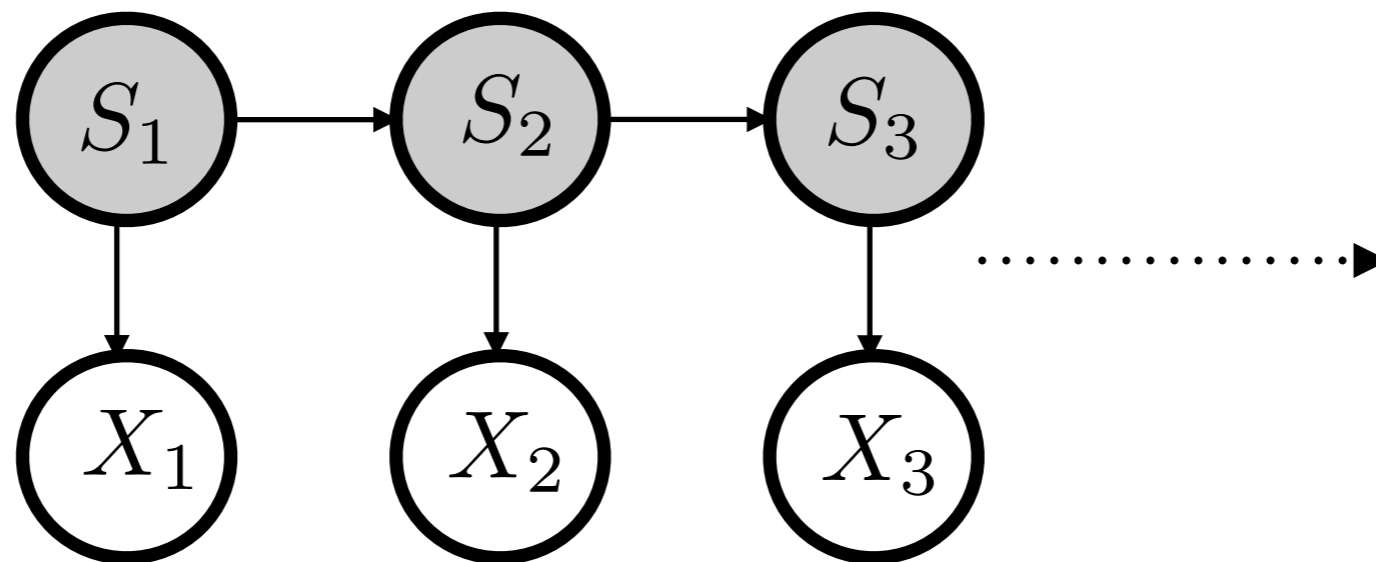
$$\text{message}_{S_{t-1} \mapsto S_t}(k) = P(S_t = k, X_1, \dots, X_{t-1})$$

$$\text{message}_{S_{t+1} \mapsto S_t}(k) = P(X_n, \dots, X_{t+1} | S_t = k)$$

Backward:

$$P(X_n, \dots, X_{t+1} | S_t = k) = \sum_{j=1}^K P(X_n, \dots, X_{t+2} | S_{t+1} = j) P(X_{t+1} | S_{t+1} = j) P(S_{t+1} = j | S_t = k)$$

# INFERENCE IN HMM



$$\text{message}_{S_{t-1} \mapsto S_t}(k) = P(S_t = k, X_1, \dots, X_{t-1})$$

$$\text{message}_{S_{t+1} \mapsto S_t}(k) = P(X_n, \dots, X_{t+1} | S_t = k)$$

Backward:

$$P(X_n, \dots, X_{t+1} | S_t = k) = \sum_{j=1}^K P(X_n, \dots, X_{t+2} | S_{t+1} = j) P(X_{t+1} | S_{t+1} = j) P(S_{t+1} = j | S_t = k)$$

$$\text{message}_{S_{t+1} \mapsto S_t}(k) = \sum_{j=1}^K \text{message}_{S_{t+2} \mapsto S_{t+1}}(j) P(X_{t+1} | S_{t+1} = j) P(S_{t+1} = j | S_t = k)$$

# LEARNING PARAMETERS FOR HMM

- Now that we have algorithm for inference, what about learning
- Given observations, how do we estimate parameters for HMM?  
Three guesses ...



# EM FOR HMM (BAUM WELCH)

- EM algorithm of course, for HMM its referred to as Baum Welch algorithm
- Initialize Transition and Emission probability tables arbitrarily
- For  $i = 1$  to convergence:

**E-step** For every state variable  $t \in \{1, \dots, n\}$ ,  
Use forward-backward algorithm to compute probabilities of latent variables given observation

**M-step** Optimize weighted log likelihood as usual:

$$\theta^{(i)} = \arg \max_{\theta \in \Theta} \sum_{S_{1,\dots,n}} P(S_{1,\dots,n} | X_{1,\dots,n}, \theta^{(i-1)}) \log P(X_{1,\dots,n}, S_{1,\dots,n} | \theta)$$

# LETS SIMPLIFY M-STEP

$$\log P(X_{1,\dots,n}, S_{1,\dots,n}|\theta) :$$

:

# LETS SIMPLIFY M-STEP

$$\log P(X_{1,\dots,n}, S_{1,\dots,n}|\theta) = \log \left( \prod_{t=1}^n P(X_t|S_t, \theta) \prod_{t=1}^n P(S_t|S_{t-1}, \theta) \right)$$

:

# LETS SIMPLIFY M-STEP

$$\begin{aligned}\log P(X_{1,\dots,n}, S_{1,\dots,n}|\theta) &= \log \left( \prod_{t=1}^n P(X_t|S_t, \theta) \prod_{t=1}^n P(S_t|S_{t-1}, \theta) \right) \\ &= \sum_{t=1}^n \log P(X_t|S_t, \theta) + \sum_{t=1}^n \log P(S_t|S_{t-1}, \theta)\end{aligned}$$

# LETS SIMPLIFY M-STEP

$$\begin{aligned}\log P(X_{1,\dots,n}, S_{1,\dots,n}|\theta) &= \log \left( \prod_{t=1}^n P(X_t|S_t, \theta) \prod_{t=1}^n P(S_t|S_{t-1}, \theta) \right) \\ &= \sum_{t=1}^n \log P(X_t|S_t, \theta) + \sum_{t=1}^n \log P(S_t|S_{t-1}, \theta)\end{aligned}$$

Hence,

$$\begin{aligned}&\sum_{S_{1,\dots,n}} P(S_{1,\dots,n}|X_{1,\dots,n}, \theta^{(i-1)}) \log P(X_{1,\dots,n}, S_{1,\dots,n}|\theta) \\ &= \sum_{t=1}^n \sum_{s_t=1}^K P(S_t = s_t|X_{1,\dots,n}, \theta^{i-1}) \log P(X_t|S_t = s_t, \theta) \\ &\quad + \sum_{t=1}^n \sum_{s_t, s_{t-1}=1}^K P(S_t = s_t, S_{t-1} = s_{t-1}|X_{1,\dots,n}, \theta^{i-1}) \log P(S_t|S_{t-1}, \theta)\end{aligned}$$

# E-STEP

# E-STEP

- Only need to compute  $P(S_t = s_t | X_{1,\dots,n}, \theta^{i-1})$  and  $P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1})$  using forward-backward

# E-STEP

- Only need to compute  $P(S_t = s_t | X_{1,\dots,n}, \theta^{i-1})$  and  $P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1})$  using forward-backward
- First term is immediate

$$P(S_t = s_t | X_{1,\dots,n}, \theta^{i-1}) \propto m_{S_{t-1} \mapsto S_t}(s_t) \cdot m_{S_{t+1} \mapsto S_t}(s_t) \cdot E^{(i-1)}[s_t, X_t]$$



# E-STEP

- Only need to compute  $P(S_t = s_t | X_{1,\dots,n}, \theta^{i-1})$  and  $P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1})$  using forward-backward
- First term is immediate

$$P(S_t = s_t | X_{1,\dots,n}, \theta^{i-1}) \propto m_{S_{t-1} \mapsto S_t}(s_t) \cdot m_{S_{t+1} \mapsto S_t}(s_t) \cdot E^{(i-1)}[s_t, X_t]$$

- For second term,

$$\begin{aligned} & P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ & \propto m_{S_{t-1} \mapsto S_t}(s_t) T^{(i-1)}[s_{t-1}, s_t] P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ & \propto m_{S_{t-1} \mapsto S_t}(s_t) T^{(i-1)}[s_{t-1}, s_t] m_{S_{t-2} \mapsto S_{t-1}}(s_{t-1}) m_{S_t \mapsto S_{t-1}}(s_{t-1}) E^{(i-1)}[s_{t-1}, X_{t-1}] \end{aligned}$$

Why?

# E-STEP

$$P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1})$$

# E-STEP

$$\begin{aligned} &P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &= P(S_t = s_t, | S_{t-1} = s_{t-1}, X_{1,\dots,n}, \theta^{t-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \end{aligned}$$

# E-STEP

$$\begin{aligned} &P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &= P(S_t = s_t, | S_{t-1} = s_{t-1}, X_{1,\dots,n}, \theta^{t-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &= P(S_t = s_t, | S_{t-1} = s_{t-1}, X_{t,\dots,n}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \end{aligned}$$

# E-STEP

$$\begin{aligned} & P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &= P(S_t = s_t, | S_{t-1} = s_{t-1}, X_{1,\dots,n}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &= P(S_t = s_t, | S_{t-1} = s_{t-1}, X_{t,\dots,n}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &\propto P(X_{t,\dots,n} | S_t = s_t, S_{t-1} = s_{t-1}, \theta^{i-1}) \\ &\quad P(S_t = s_t | S_{t-1} = s_{t-1}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \end{aligned}$$

# E-STEP

$$\begin{aligned} & P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &= P(S_t = s_t, | S_{t-1} = s_{t-1}, X_{1,\dots,n}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &= P(S_t = s_t, | S_{t-1} = s_{t-1}, X_{t,\dots,n}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &\propto P(X_{t,\dots,n} | S_t = s_t, S_{t-1} = s_{t-1}, \theta^{i-1}) \\ &\quad P(S_t = s_t | S_{t-1} = s_{t-1}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &\propto P(X_{t,\dots,n} | S_t = s_t, \theta^{i-1}) \\ &\quad T^{(i-1)}[s_{t-1}, s_t] P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \end{aligned}$$

# E-STEP

$$\begin{aligned} & P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &= P(S_t = s_t, | S_{t-1} = s_{t-1}, X_{1,\dots,n}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &= P(S_t = s_t, | S_{t-1} = s_{t-1}, X_{t,\dots,n}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &\propto P(X_{t,\dots,n} | S_t = s_t, S_{t-1} = s_{t-1}, \theta^{i-1}) \\ &\quad P(S_t = s_t | S_{t-1} = s_{t-1}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &\propto P(X_{t,\dots,n} | S_t = s_t, \theta^{i-1}) \\ &\quad T^{(i-1)}[s_{t-1}, s_t] P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &\propto m_{S_{t-1} \mapsto S_t}(s_t) \cdot T^{(i-1)}[s_{t-1}, s_t] \cdot P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \end{aligned}$$

# E-STEP

$$\begin{aligned} & P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &= P(S_t = s_t, | S_{t-1} = s_{t-1}, X_{1,\dots,n}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &= P(S_t = s_t, | S_{t-1} = s_{t-1}, X_{t,\dots,n}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &\propto P(X_{t,\dots,n} | S_t = s_t, S_{t-1} = s_{t-1}, \theta^{i-1}) \\ &\quad P(S_t = s_t | S_{t-1} = s_{t-1}, \theta^{i-1}) P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &\propto P(X_{t,\dots,n} | S_t = s_t, \theta^{i-1}) \\ &\quad T^{(i-1)}[s_{t-1}, s_t] P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &\propto m_{S_{t-1} \mapsto S_t}(s_t) \cdot T^{(i-1)}[s_{t-1}, s_t] \cdot P(S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1}) \\ &\propto m_{S_{t-1} \mapsto S_t}(s_t) \cdot T^{(i-1)}[s_{t-1}, s_t] \\ &\quad m_{S_{t-2} \mapsto S_{t-1}}(s_{t-1}) \cdot m_{S_t \mapsto S_{t-1}}(s_{t-1}) \cdot E^{(i-1)}[s_{t-1}, X_{t-1}] \end{aligned}$$



# BAUM WELCH ALGORITHM

# BAUM WELCH ALGORITHM

Initialize  $T^0, E^0$  probability tables

# BAUM WELCH ALGORITHM

Initialize  $T^0, E^0$  probability tables

For  $i = 1$  to convergence

# BAUM WELCH ALGORITHM

Initialize  $T^0, E^0$  probability tables

For  $i = 1$  to convergence

- E-step:
  - Run Forward-Backward algorithm and compute messages

# BAUM WELCH ALGORITHM

Initialize  $T^0, E^0$  probability tables

For  $i = 1$  to convergence

- E-step:
  - Run Forward-Backward algorithm and compute messages
  - For every  $t$  compute  $P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1})$  and  $P(S_t = s_t | X_{1,\dots,n}, \theta^{i-1})$  as in previous slides

# BAUM WELCH ALGORITHM

Initialize  $T^0, E^0$  probability tables

For  $i = 1$  to convergence

- E-step:
  - Run Forward-Backward algorithm and compute messages
  - For every  $t$  compute  $P(S_t = s_t, S_{t-1} = s_{t-1} | X_{1,\dots,n}, \theta^{i-1})$  and  $P(S_t = s_t | X_{1,\dots,n}, \theta^{i-1})$  as in previous slides
- M-step:

$$\forall u, v \quad T^{(i)}[u, v] = \frac{\sum_{t=2}^n P(S_t = v, S_{t-1} = u | X_{1,\dots,n}, \theta^{i-1})}{\sum_{t=2}^n P(S_{t-1} = u | X_{1,\dots,n}, \theta^{i-1})}$$

$$\forall v, e \quad E^{(i)}[v, e] = \frac{\sum_{t=1}^n P(S_t = v | X_{1,\dots,n}, \theta^{i-1}) \cdot \mathbf{1}_{X_t=e}}{\sum_{t=1}^n P(S_t = v | X_{1,\dots,n}, \theta^{i-1})}$$