

# COMS W4701: Artificial Intelligence

## Lecture 17: Hidden Markov Models

Tony Dear, Ph.D.

Department of Computer Science

School of Engineering and Applied Sciences

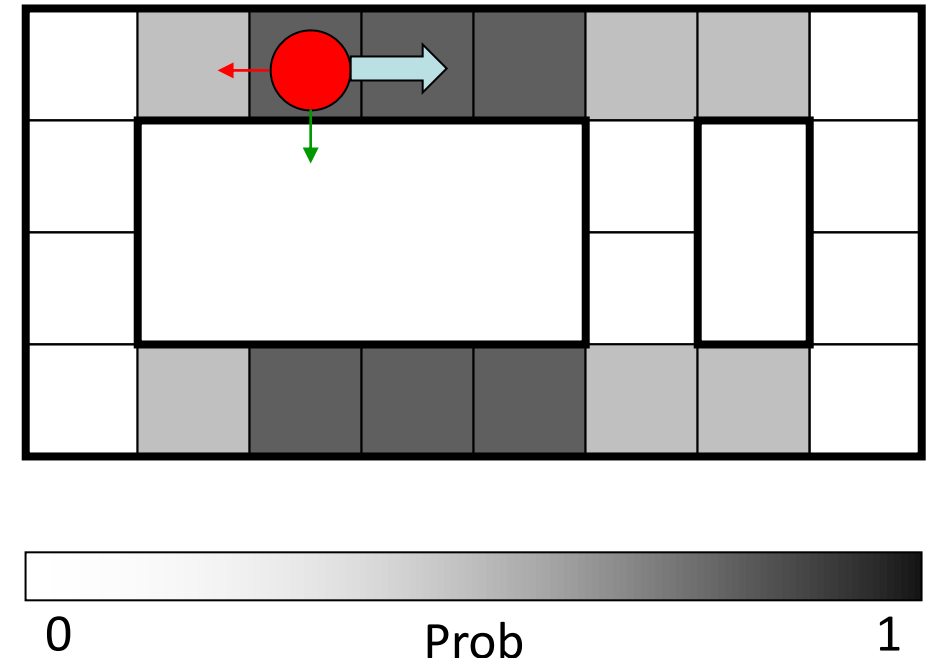
# Today

---

- Markov chains
- Hidden Markov models
- State estimation (filtering): Forward algorithm

# Temporal Reasoning

- Scenario: An agent's state changes over time, but not directly observable
- *Belief state*: A random variable  $X_t$  representing the agent's current state, along with a probability distribution over the state space
- A probabilistic *transition model* describes how  $X_t$  is derived from past states
- We will be interested in looking at how  $X_t$  changes over time, possibly incorporating sensor information



# Markov Chains

- **Markov chain:** A sequence of RVs  $X_1, X_2, \dots$ , s.t.  $X_t$  only depends on  $X_{t-1}$
- Parameters: Initial state  $P(X_1)$ , **transition model**  $P(X_t|X_{t-1})$

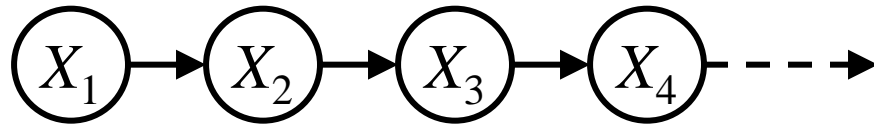
- If  $|X_t| = n$ , we have  $n^2$  different  $P(x_t|x_{t-1})$  transition probabilities
- Define a  $n \times n$  *transition matrix*  $T$ , where  $T_{ij} = P(X_t = j \mid X_{t-1} = i)$

$$T = \begin{bmatrix} P(X_t = 1 \mid X_{t-1} = 1) & \cdots & P(X_t = n \mid X_{t-1} = 1) \\ \vdots & \ddots & \vdots \\ P(X_t = 1 \mid X_{t-1} = n) & \cdots & P(X_t = n \mid X_{t-1} = n) \end{bmatrix}$$

- Sum of each row  $\sum_j T_{ij} = \sum_j P(X_t = j \mid X_{t-1} = i) = 1$

# Markov Assumption

- **Markov assumption:**  $X_t$  is independent of all past states given  $X_{t-1}$



$$X_t \perp\!\!\!\perp X_1, \dots, X_{t-2} \mid X_{t-1}$$

$$X_3 \perp\!\!\!\perp X_1 \mid X_2$$

$$X_4 \perp\!\!\!\perp X_1, X_2 \mid X_3$$

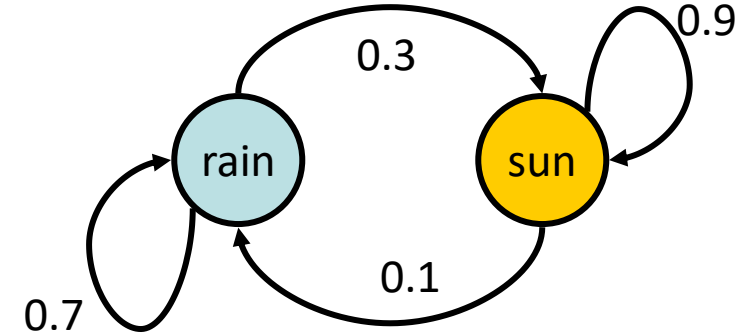
- Chain rule for joint distribution can be greatly simplified!

$$P(X_1, X_2, \dots, X_T) = P(X_1)P(X_2|X_1)P(X_3|X_2) \dots P(X_T|X_{T-1})$$

$$= P(X_1) \prod_{t=2}^T P(X_t|X_{t-1})$$

# Example: Markov Chains

$$P(X_1) = \begin{matrix} & \text{rain} & \text{sun} \\ \begin{matrix} \text{rain} \\ \text{sun} \end{matrix} & \begin{pmatrix} 0.8 & 0.2 \end{pmatrix} \end{matrix} \quad T = \begin{matrix} & \text{rain} & \text{sun} \\ \begin{matrix} \text{rain} \\ \text{sun} \end{matrix} & \begin{pmatrix} 0.7 & 0.3 \\ 0.1 & 0.9 \end{pmatrix} \end{matrix}$$



- $P(X_2 = \text{rain}) = \sum_{x_1} P(x_1)P(X_2 = \text{rain}|x_1) = 0.8(0.7) + 0.2(0.1) = 0.58$
- $P(X_2 = \text{sun}) = \sum_{x_1} P(x_1)P(X_2 = \text{sun}|x_1) = 0.8(0.3) + 0.2(0.9) = 0.42$
- Alternatively, can compute  $P(X_2) = P(X_1)T$ ,  $P(X_3) = P(X_2)T$ , ...,  $P(X_t) = P(X_{t-1})T$
- More generally,  $P(X_t) = P(X_1)T^{t-1}$

# Stationary Distributions

- Observation:  $\pi = (.25 \ .75)$  satisfies  $\pi = \pi \cdot T$
- $\pi$  is an *eigenvector* of  $T^\top$  corresponding to eigenvalue 1
- $\pi$  is a **stationary distribution** of this transition matrix  $T = \begin{pmatrix} 0.7 & 0.3 \\ 0.1 & 0.9 \end{pmatrix}$
- All transition matrices have at least one stationary distribution
- Find the appropriate *eigenvector*  $\pi$  of  $T^\top$  and rescale as  $\pi / \sum_i \pi_i$  to ensure that the vector sum is 1
- Some Markov chains may have multiple stationary distributions

# Markov Chain Applications

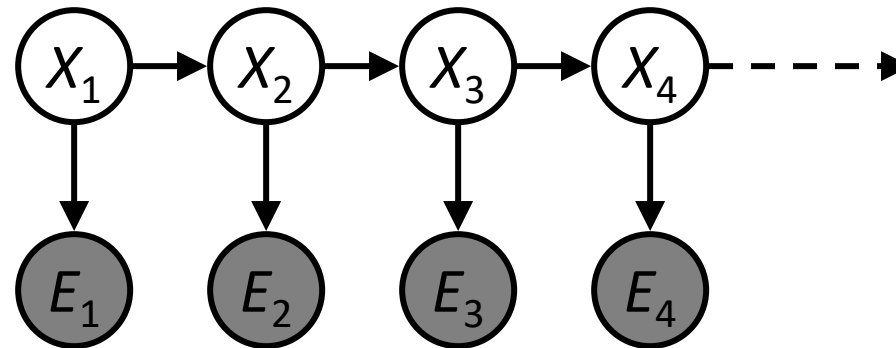
---

- Bioinformatics, population dynamics, epidemic modeling
- Thermodynamics, statistical mechanics, chemical reaction modeling
- Queuing theory, income and market modeling, game modeling
- Speech recognition and text generation, n-gram models
  - Unigram model:  $P(word_t = i)$ , bigram model:  $P(word_t = i \mid word_{t-1} = j)$
- Web browsing: PageRank algorithm to determine webpage traffic
  - Model probabilities of navigating to existing outgoing link or arbitrary webpage

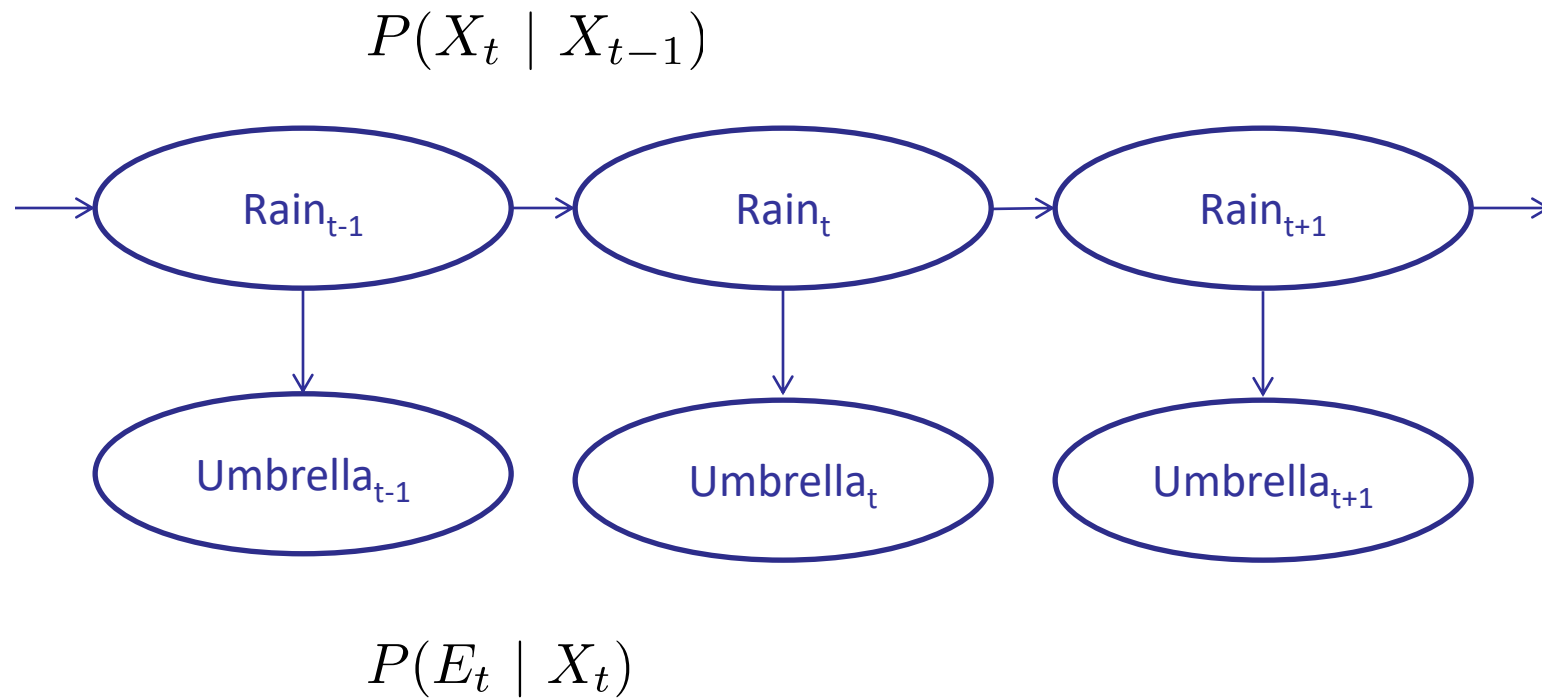


# Hidden Markov Models

- With Markov chains, we do not directly observe the state
- Now let's suppose we can observe indirect *evidence* of states
- **Hidden Markov model:** A Markov process with *hidden* states  $X_t$  and *observable* evidence variables  $E_t$
- Initial belief state:  $P(X_1)$
- Transition model:  $P(X_t|X_{t-1})$
- Observation model:  $P(E_t|X_t)$



# Example: Weather HMM



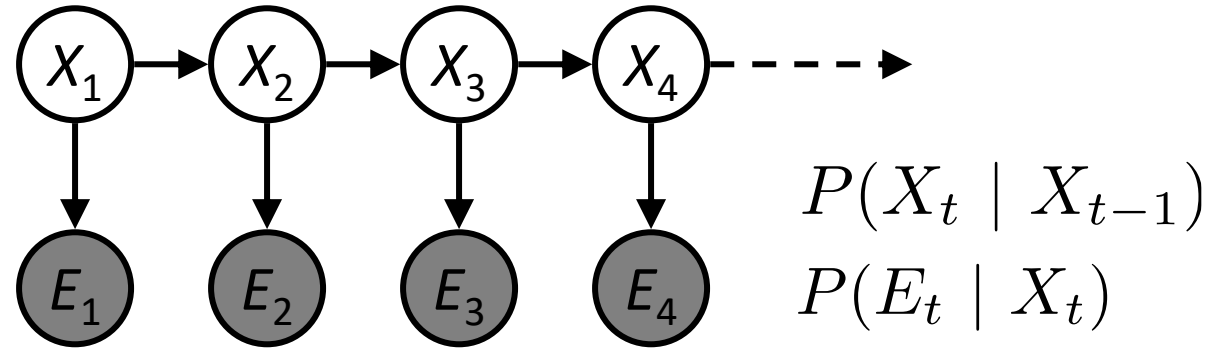
$X_{t-1}$	$X_t$	$P(X_t X_{t-1})$
+r	+r	0.7
+r	-r	0.3
-r	+r	0.3
-r	-r	0.7

$X_t$	$E_t$	$P(E_t X_t)$
+r	+u	0.9
+r	-u	0.1
-r	+u	0.2
-r	-u	0.8

# Conditional Independences

- Markov chain independences:

$$X_t \perp\!\!\!\perp X_1, \dots, X_{t-2} \mid X_{t-1}$$



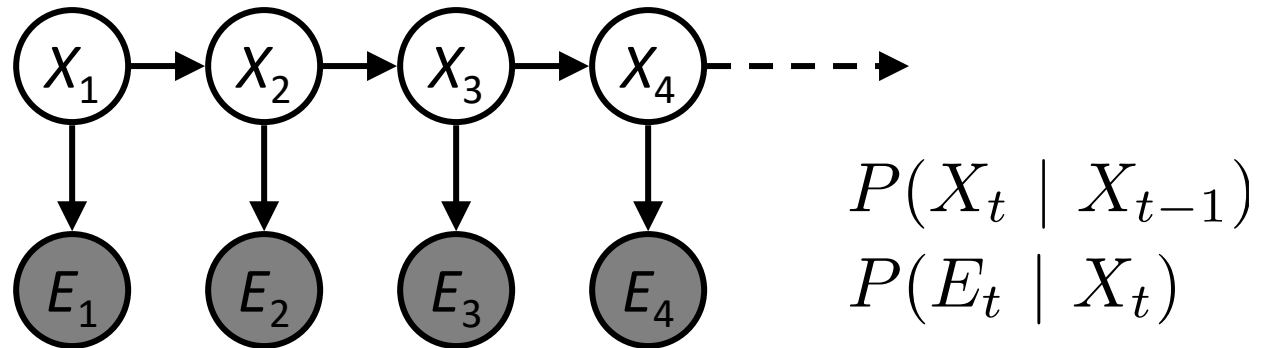
- A state is conditionally independent of past states and evidence given preceding state:

$$X_t \perp\!\!\!\perp X_1, E_1, \dots, X_{t-2}, E_{t-2}, E_{t-1} \mid X_{t-1}$$

- An observation is conditionally independent of past states and evidence given current state:

$$E_t \perp\!\!\!\perp X_1, E_1, \dots, X_{t-2}, E_{t-2}, X_{t-1}, E_{t-1} \mid X_t$$

# Joint Distribution



- General joint distribution:

$$P(X_1, E_1, \dots, X_T, E_T) = P(X_1)P(E_1|X_1) \prod_{t=2}^T P(X_t|X_{t-1})P(E_t|X_t)$$

- Marginal or smaller joint distributions can be found by summing out RVs
- For certain computations we don't even need the entire joint distribution!

# Inference

---

- Inference tasks compute belief states or hidden states given evidence
- **Filtering (state estimation):** Find  $P(X_t \mid e_{1:t})$ 
  - Estimate the belief state, given a sequence of past observations
- **Decoding:** Find  $\operatorname{argmax}_{x_{1:t}} P(x_{1:t} \mid e_{1:t})$ 
  - Find the *sequence* of hidden states that best explains given observations
- **Smoothing:** Find  $P(X_k \mid e_{1:t})$ , for  $1 \leq k < t$ 
  - Use both past and future evidence to *smooth* a belief state

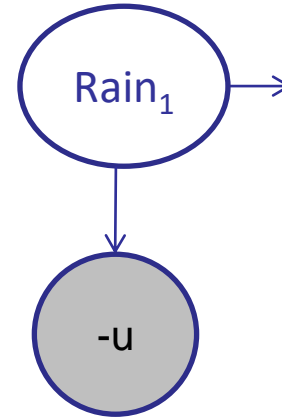
# Example: Weather HMM

- Want to find  $P(X_1|e_1) = \frac{P(X_1, e_1)}{P(e_1)}$
- $P(e_1)$  is a *constant* for all values of  $X_1$  since  $e_1$  is already observed (fixed)!

$$P(X_1|e_1) \propto P(X_1, e_1) = P(e_1|X_1) * P(X_1)$$

$$= (0.1 \quad 0.8) * (0.5 \quad 0.5) = (0.05 \quad 0.4)$$

- Since we know  $P(X_1|e_1)$  sums to 1,  $P(e_1)$  must be equal to  $\sum_{x_1} P(x_1, e_1)$



$$P(X_1) = (0.5 \quad 0.5)$$

$X_t$	$E_t$	$P(E_t X_t)$
+r	+u	0.9
+r	-u	0.1
-r	+u	0.2
-r	-u	0.8

$$P(X_1|e_1) = \frac{P(X_1, e_1)}{P(e_1)} = \frac{(0.05 \quad 0.4)}{0.05 + 0.4}$$

$$= (0.11 \quad 0.89)$$

# State Estimation

---

- In a Markov chain, we obtained  $P(X_{t+1})$  from  $P(X_t)$  by multiplication with transition probabilities
- We just showed how to obtain  $P(X_{t+1}|e_{t+1})$  from  $P(X_{t+1})$  by multiplication with observation probabilities, followed by normalization
- To efficiently solve the state estimation problem of finding  $P(X_{t+1}|e_{1:t+1})$ , we need to show how to perform these steps starting from  $P(X_t|e_{1:t})$
- (To simplify calculations, we will work primarily with  $P(X_t, e_{1:t})$ )

# Forward Algorithm

- Given  $P(x_t, e_{1:t})$ : Conditional independence

$$\sum_{x_t} P(X_{t+1} | x_t, e_{1:t}) P(\mathbf{x}_t, \mathbf{e}_{1:t}) = \sum_{x_t} P(X_{t+1}, x_t, e_{1:t}) = P(X_{t+1}, \mathbf{e}_{1:t})$$

- So we have  $P(X_{t+1}, e_{1:t}) = P(X_t, e_{1:t}) \cdot T$  (same as Markov chains)

- Product rule:  $P(e_{t+1} | X_{t+1}, e_{1:t}) P(X_{t+1}, \mathbf{e}_{1:t}) = P(X_{t+1}, \mathbf{e}_{1:t+1})$   
Conditional independence

- So  $P(X_{t+1}, e_{1:t+1}) = P(X_{t+1}, e_{1:t}) * O_{t+1}$ , where  $O_{t+1} = P(e_{t+1} | X_{t+1})$  is a vector of observation probabilities and  $*$  is an elementwise product



# Forward Algorithm

- Given:  $\alpha_0 = P(X_0)$ , or if starting with  $\alpha_1 = P(X_1)$ , skip the first “elapse time” step and start by observing evidence  $e_1$

$$\alpha_t = P(X_t, e_{1:t})$$

- For each timestep  $t$ :

$$\alpha'_{t+1} = P(X_{t+1}, e_{1:t})$$

- Elapse time:

$$\alpha'_{t+1} = \alpha_t T$$

$$\alpha_{t+1} = P(X_{t+1}, e_{1:t+1})$$

- Observe evidence  $e_{t+1}$ :

$$\alpha_{t+1} = \alpha'_{t+1} * O_{t+1}$$

- Normalize (as needed):

$$P(X_{t+1} | e_{1:t+1}) = \alpha_{t+1} / \sum \alpha_{t+1}$$

# Example: Weather HMM

$$T = \begin{pmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{pmatrix} \begin{matrix} +r \\ -r \end{matrix}$$

Suppose  $\alpha_0 = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$

$$\alpha'_1 = \alpha_0^\top T = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$$

$$\alpha'_2 = \alpha_1^\top T = \begin{pmatrix} .155 \\ .295 \end{pmatrix}$$

$$\alpha'_3 = \alpha_2^\top T = \begin{pmatrix} .115 \\ .083 \end{pmatrix}$$

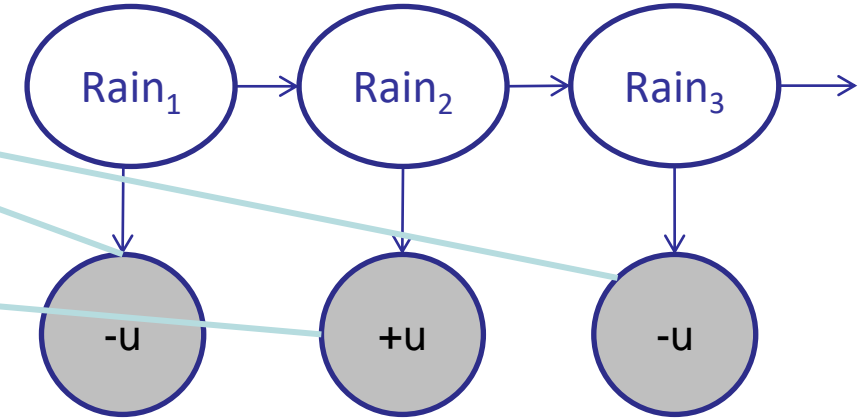
$$O_1 = O_3 = \begin{pmatrix} 0.1 \\ 0.8 \end{pmatrix}$$

$$O_2 = \begin{pmatrix} 0.9 \\ 0.2 \end{pmatrix}$$

$$\alpha_1 = \alpha'_1 * O_1 = \begin{pmatrix} 0.05 \\ 0.4 \end{pmatrix}$$

$$\alpha_2 = \alpha'_2 * O_2 = \begin{pmatrix} .1395 \\ .059 \end{pmatrix}$$

$$\alpha_3 = \alpha'_3 * O_3 = \begin{pmatrix} .0115 \\ .0664 \end{pmatrix}$$



$$P(X_1|e_1) = \begin{pmatrix} 0.11 \\ 0.89 \end{pmatrix}$$

$$P(X_2|e_{1:2}) = \begin{pmatrix} .703 \\ .297 \end{pmatrix}$$

$$P(X_3|e_{1:3}) = \begin{pmatrix} .148 \\ .852 \end{pmatrix}$$

# Summary

---

- Temporal models are used to track partially observable environments
- Maintain and update belief states (probability distributions)
- Markov chains may have stationary distributions or steady state behavior
- Inference in HMMs compute hidden information given observed information
- State estimation: Forward algorithm iteratively computes the current state distribution given evidence to date