Lecture 23: 17 April, 2025

Madhavan Mukund https://www.cmi.ac.in/~madhavan

Data Mining and Machine Learning January–April 2025

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

D-Separation

Conditional Independence

• Check if $X \perp Y \mid Z$

- Dependence should be blocked on every trail from X to Y
 - Each undirected path from X to Y is a sequence of basic trails
 - For (a), (b), (c), need Z present
 - For (d), need Z absent
 - In general, V-structure includes descendants of the bottom node
- x and y are D-separated given z if all trails are blocked
- Variation of breadth first search (BFS) to check if y is reachable from x through some trail
- Extends to sets each $x \in X$ is D-separated from each $y \in Y$



Graph properh

■ *MB*(*X*) — Markov blanket of *X*



3/16

• MB(X) — Markov blanket of X





- *MB*(*X*) Markov blanket of *X*
 - Parents(X)
 - Children(X)



3/16

- *MB*(*X*) Markov blanket of *X*
 - Parents(X)
 - Children(X)
 - Parents of Children(X)



3/16

- *MB*(*X*) Markov blanket of *X*
 - \blacksquare Parents(X)
 - Children(X)
 - Parents of Children(X)
- $X \perp \neg MB(X) \mid MB(X)$



John and Mary call Pearl. What is the probability that there has been a burglary?



▶ < ∃ ▶</p>

- John and Mary call Pearl. What is the probability that there has been a burglary?
- Want $P(b \mid m, j)$



▶ < ∃ ▶</p>

- John and Mary call Pearl. What is the probability that there has been a burglary?
- Want $P(b \mid m, j)$

 $\bullet \frac{P(b,m,j)}{P(m,j)}$



< ∃ >

- John and Mary call Pearl. What is the probability that there has been a burglary?
- Want $P(b \mid m, j)$
- $\blacksquare \frac{P(b,m,j)}{P(m,j)}$
- Use chain rule to evaluate joint probabilities



< E

- John and Mary call Pearl. What is the probability that there has been a burglary?
- Want $P(b \mid m, j)$
- $\blacksquare \frac{P(b,m,j)}{P(m,j)}$
- Use chain rule to evaluate joint probabilities
- Reorder variables appropriately, topological order of graph



•
$$P(m,j,b) = P(b) \sum_{e=0}^{1} P(e) \sum_{a=0}^{1} P(a \mid b, e) P(m \mid a) P(j \mid a)$$

Image: A image: A

•
$$P(m,j,b) = P(b) \sum_{e=0}^{1} P(e) \sum_{a=0}^{1} P(a \mid b, e) P(m \mid a) P(j \mid a)$$



•
$$P(m,j,b) = P(b) \sum_{e=0}^{1} P(e) \sum_{a=0}^{1} P(a \mid b, e) P(m \mid a) P(j \mid a)$$

- Construct the computation tree
- Use dynamic programming to avoid duplicated computations



•
$$P(m,j,b) = P(b) \sum_{e=0}^{1} P(e) \sum_{a=0}^{1} P(a \mid b, e) P(m \mid a) P(j \mid a)$$

- Construct the computation tree
- Use dynamic programming to avoid duplicated computations
- However, exact inference is NP-complete, in general





 Instead, approximate inference through sampling

.70

.01

.01

.70

Generate random samples
 (b, e, a, m, j), count to estimate probabilities



< □ > < 円

- Generate random samples
 (b, e, a, m, j), count to estimate probabilities
- Random samples should respect conditional probabilities



→ < ∃→

- Generate random samples
 (b, e, a, m, j), count to estimate probabilities
- Random samples should respect conditional probabilities
- Fix parents of x before generating x

P(6/mj)



< ∃ >

- Generate random samples
 (b, e, a, m, j), count to estimate probabilities
- Random samples should respect conditional probabilities
- Fix parents of x before generating x
- Generate in topological order
 - Generate b, e with probabilities P(b) and P(e)
 - Generate *a* with probability *P*(*a* | *b*, *e*)
 - Generate *j*, *m* with probabilities *P*(*j* | *a*), *P*(*m* | *a*)



• We are interested in $P(b \mid j, m)$



3



3

- We are interested in $P(b \mid j, m)$
- Samples with $\neg j$ or $\neg m$ are useless
- Can we sample more efficiently?



→ < ∃→

3

■ *P*(*Rain* | *Cloudy*, *Wet Grass*)



▶ ▲ 国 ▶ ▲ 国 ▶

< □ > < 円

- *P*(*Rain* | *Cloudy*, *Wet Grass*)
- Topological order
 - Generate *Cloudy*
 - Generate Sprinkler, Rain
 - Generate Wet Grass



< ∃ >

- *P*(*Rain* | *Cloudy*, *Wet Grass*)
- Topological order
 - Generate *Cloudy*
 - Generate Sprinkler, Rain
 - Generate Wet Grass
- If we start with ¬*Cloudy*, sample is useless



< 3

э

8/16

- *P*(*Rain* | *Cloudy*, *Wet Grass*)
- Topological order
 - Generate Cloudy
 - Generate Sprinkler, Rain
 - Generate Wet Grass
- If we start with ¬*Cloudy*, sample is useless
- Immediately stop and reject this sample — rejection sampling



- *P*(*Rain* | *Cloudy*, *Wet Grass*)
- Topological order
 - Generate Cloudy
 - Generate Sprinkler, Rain
 - Generate Wet Grass
- If we start with ¬*Cloudy*, sample is useless
- Immediately stop and reject this sample rejection sampling
- General problem with low probability situation — many samples are rejected



■ *P*(*Rain* | *Cloudy*, *Wet Grass*)



• • = •

■ *P*(*Rain* | *Cloudy*, *Wet Grass*)

■ Fix evidence *Cloudy*, *Wet Grass* true



▶ < ∃ ▶</p>

- *P*(*Rain* | *Cloudy*, *Wet Grass*)
- Fix evidence Cloudy, Wet Grass true
- Then generate the other variables



< ∃

- *P*(*Rain* | *Cloudy*, *Wet Grass*)
- Fix evidence *Cloudy*, *Wet Grass* true
- Then generate the other variables
- Suppose we generate $c, \neg s, r, w$



- P(Rain | Cloudy, Wet Grass)
- Fix evidence *Cloudy*, *Wet Grass* true
- Then generate the other variables
- Suppose we generate $c, \neg s, r, w$
- Compute likelihood of evidence: 0.5 × 0.9 = 0.45

$$P(c) \times P(w) \quad \text{given this comple}$$

$$0.5 \times 0.9 = 0.45$$



- P(Rain | Cloudy, Wet Grass)
- Fix evidence *Cloudy*, *Wet Grass* true
- Then generate the other variables
- Suppose we generate $c, \neg s, r, w$
- Compute likelihood of evidence: 0.5 × 0.9 = 0.45
- 0.45 is likelihood weight of sample



- P(Rain | Cloudy, Wet Grass)
- Fix evidence *Cloudy*, *Wet Grass* true
- Then generate the other variables
- Suppose we generate $c, \neg s, r, w$
- Compute likelihood of evidence: 0.5 × 0.9 = 0.45
- 0.45 is likelihood weight of sample
- Samples *s*₁, *s*₂, ..., *s*_N with weights *w*₁, *w*₂, ... *w*_N


Likelihood weighted sampling

- P(Rain | Cloudy, Wet Grass)
- Fix evidence *Cloudy*, *Wet Grass* true
- Then generate the other variables
- Suppose we generate $c, \neg s, r, w$
- Compute likelihood of evidence: 0.5 × 0.9 = 0.45
- 0.45 is likelihood weight of sample
 Samples s₁, s₂, ..., s_N with weights w₁, w₂, ... w_N

•
$$P(r \mid c, w) = \frac{\sum_{s_i \text{ has rain } W_i}}{\sum_{1 \le j \le N} w_j}$$



• State of a Bayesian network is a valuation of variables (V_1, V_2, \ldots, V_n)

► < ∃ ►</p>

3

- State of a Bayesian network is a valuation of variables (V_1, V_2, \ldots, V_n)
- Move probabilistically from $s_j = (x_1, x_2, \dots, x_n)$ to $s_k = (y_1, y_2, \dots, y_n)$

▶ < ∃ ▶</p>

- State of a Bayesian network is a valuation of variables (V_1, V_2, \ldots, V_n)
- Move probabilistically from $s_j = (x_1, x_2, \dots, x_n)$ to $s_k = (y_1, y_2, \dots, y_n)$
- Allow such a move only when s_j , s_k differ at exactly one position
 - $s_j = (x_1, x_2, ..., x_{i-1}, x_i, x_{i+1}, ..., x_n)$
 - $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$

- State of a Bayesian network is a valuation of variables (V_1, V_2, \ldots, V_n)
- Move probabilistically from $s_j = (x_1, x_2, \dots, x_n)$ to $s_k = (y_1, y_2, \dots, y_n)$
- Allow such a move only when s_j , s_k differ at exactly one position
 - $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$
 - $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$
- Sampling algorithm
 - Current state is $s_j = (x_1, x_2, \dots, x_n)$
 - Choose *i* uniformly in [1, *n*]
 - Resample x_i given current values $(x_1, x_2, \ldots, x_{i-1}, x_{i+1}, \ldots, x_n)$
 - Random walk through state space count number of visits to each state

3

- State of a Bayesian network is a valuation of variables (V_1, V_2, \ldots, V_n)
- Move probabilistically from $s_j = (x_1, x_2, \dots, x_n)$ to $s_k = (y_1, y_2, \dots, y_n)$
- Allow such a move only when s_j , s_k differ at exactly one position
 - $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$
 - $s_k = (x_1, x_2, \ldots, x_{i-1}, y_i, x_{i+1}, \ldots, x_n)$
- Sampling algorithm
 - Current state is $s_j = (x_1, x_2, \dots, x_n)$
 - Choose *i* uniformly in [1, *n*]
 - Resample x_i given current values $(x_1, x_2, \ldots, x_{i-1}, x_{i+1}, \ldots, x_n)$
 - Random walk through state space count number of visits to each state
- Need to compute $P[y_i | x_1, x_2, ..., x_{i-1}, x_{i+1}, ..., x_n]$

- State of a Bayesian network is a valuation of variables (V_1, V_2, \ldots, V_n)
- Move probabilistically from $s_j = (x_1, x_2, \dots, x_n)$ to $s_k = (y_1, y_2, \dots, y_n)$
- Allow such a move only when s_j , s_k differ at exactly one position
 - $s_j = (x_1, x_2, \ldots, x_{i-1}, x_i, x_{i+1}, \ldots, x_n)$
 - $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$
- Sampling algorithm
 - Current state is $s_j = (x_1, x_2, \dots, x_n)$
 - Choose *i* uniformly in [1, *n*]
 - Resample x_i given current values $(x_1, x_2, \ldots, x_{i-1}, x_{i+1}, \ldots, x_n)$
 - Random walk through state space count number of visits to each state
- Need to compute $P[y_i | x_1, x_2, ..., x_{i-1}, x_{i+1}, ..., x_n]$
- Why does this work at all?

Approximate inference using Markov chains

Markov chains

Finite set of states, with transition probabilities between states





- Finite set of states, with transition probabilities between states
- For us, a state will be an assignment of values to variables

- Finite set of states, with transition probabilities between states
- For us, a state will be an assignment of values to variables
- A three state Markov Chain



- Finite set of states, with transition probabilities between states
- For us, a state will be an assignment of values to variables
- A three state Markov Chain





- Finite set of states, with transition probabilities between states
- For us, a state will be an assignment of values to variables
- A three state Markov Chain





- Finite set of states, with transition probabilities between states
- For us, a state will be an assignment of values to variables
- A three state Markov Chain



Represent using a transition matrix — stochastic

$$A = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix}$$

P[*j*] is probability of being in state *j*

- Finite set of states, with transition probabilities between states
- For us, a state will be an assignment of values to variables
- A three state Markov Chain



Represent using a transition matrix — stochastic
$$A = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix}$$

$$P[j] \text{ is probability of being in state } j$$
Start in state 1, so initially $P = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ ($z = b \in C$)
$$A = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

$$A$$

Markov chains . . .

After one step:

$$P^{\top}A = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$



A 国
 A 国
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

< □ > < 同

3

Markov chains . . .

After one step:

$$P^{\top}A = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

• After second step: $\begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} = \begin{bmatrix} \frac{3}{4} & \frac{1}{4} & 0 \end{bmatrix}$



→

Markov chains . . .

After one step:

$$P^{\top}A = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

• After second step: $\begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} = \begin{bmatrix} \frac{3}{4} & \frac{1}{4} & 0 \end{bmatrix}$

 After k steps, P[j] is probability of being in state j



Markov chains

- After one step: $P^{\top}A = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$
- After second step: $\begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 \\ 1 & 1 & 0 \end{bmatrix} = \begin{bmatrix} \frac{3}{4} & \frac{1}{4} & 0 \end{bmatrix}$
- After k steps, P[i] is probability of being in state *j*

 $\rightarrow \begin{bmatrix} \frac{1}{4} & \frac{3}{8} & \frac{3}{8} \end{bmatrix} \rightarrow \begin{bmatrix} \frac{9}{16} \end{bmatrix}$

Continuing our example,

Madhavan Mukund

Lecture 23: 17 April. 2025

LP. PL P2]-A = [P.P.B] $\frac{1}{2}$ 2



Is it the case that P[j] > 0 for all j continuously, after some point?



Image: A (1) and A (1)

- Is it the case that P[j] > 0 for all j continuously, after some point?
- Markov chain A is ergodic if there is some t₀ such that for every P, for all t > t₀, for every j, (P^TA^t)[j] > 0.



- Is it the case that P[j] > 0 for all j continuously, after some point?
- Markov chain A is ergodic if there is some t₀ such that for every P, for all t > t₀, for every j, (P^TA^t)[j] > 0.
 - No matter where we start, after t > t₀ steps, every state has a nonzero probability of being visited in step t



- Is it the case that P[j] > 0 for all j continuously, after some point?
- Markov chain A is ergodic if there is some t₀ such that for every P, for all t > t₀, for every j, (P^TA^t)[j] > 0.
 - No matter where we start, after t > t₀ steps, every state has a nonzero probability of being visited in step t
- Properties of ergodic Markov chains
 - There is a stationary distribution π , $\pi^{\top} A = \pi$
 - **\pi** is a left eigenvector of *A*



- Is it the case that P[j] > 0 for all j continuously, after some point?
- Markov chain A is ergodic if there is some t₀ such that for every P, for all t > t₀, for every j, (P^TA^t)[j] > 0.
 - No matter where we start, after t > t₀ steps, every state has a nonzero probability of being visited in step t
- Properties of ergodic Markov chains
 - There is a stationary distribution π , $\pi^{\top} A = \pi$
 - π is a left eigenvector of A
 - For any starting distribution P, $\lim_{t\to\infty} P^{\top} A^t = \pi$







Ergodicity . . .

• How can ergodicity fail?



DMML Jan–Apr 2025

A 国
 A 国
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

< □ > < 同

3

Ergodicity . . .

- How can ergodicity fail?
 - Starting from *i*, we reach a set of states from which there is no path back to *i*



→ < ∃→

Ergodicity ...

- How can ergodicity fail?
 - Starting from *i*, we reach a set of states from which there is no path back to *i*
 - We have a cycle i → j → k → i → j → k ···, so we can only visit some states periodically



Ergodicity ...

- How can ergodicity fail?
 - Starting from *i*, we reach a set of states from which there is no path back to *i*
 - We have a cycle i → j → k → i → j → k ···, so we can only visit some states periodically
- Sufficient conditions for ergodicity



Ergodicity . . .

- How can ergodicity fail?
 - Starting from *i*, we reach a set of states from which there is no path back to *i*
 - We have a cycle i → j → k → i → j → k ···, so we can only visit some states periodically
- Sufficient conditions for ergodicity
 - Irreducibility: When viewed as a directed graph, A is strongly connected
 - For all states i, j, there is a path from i to j and a path from j to i



Ergodicity . . .

- How can ergodicity fail?
 - Starting from *i*, we reach a set of states from which there is no path back to *i*
 - We have a cycle i → j → k → i → j → k ···, so we can only visit some states periodically
- Sufficient conditions for ergodicity
 - Irreducibility: When viewed as a directed graph, A is strongly connected
 - For all states i, j, there is a path from i to j and a path from j to i
 - Aperiodicity: For any pair of vertices *i*, *j*, the gcd of the lengths of all paths from *i* to *j* is 1
 - In particular, paths (loops) from *i* to *i* do not all have lengths that are multiples of some k ≥ 2 prevents bad cycles



- Can efficiently approximate $\lim_{t\to\infty} P^{\top} A^t$ by repeated squaring: $P^{\top} A^2$, $P^{\top} A^4$, $P^{\top} A^8$, ..., $P^{\top} A^{2^k}$, ...
 - Mixing time how fast this converges to π



→ < ∃→

- Can efficiently approximate $\lim_{t\to\infty} P^{\top} A^t$ by repeated squaring: $P^{\top} A^2$, $P^{\top} A^4$, $P^{\top} A^8$, ..., $P^{\top} A^{2^k}$, ...
 - Mixing time how fast this converges to π
- Stationary distribution represents fraction of visits to each state in a long enough execution



- Can efficiently approximate $\lim_{t\to\infty} P^{\top} A^t$ by repeated squaring: $P^{\top} A^2$, $P^{\top} A^4$, $P^{\top} A^8$, ..., $P^{\top} A^{2^k}$, ...
 - Mixing time how fast this converges to π
- Stationary distribution represents fraction of visits to each state in a long enough execution
- Can we create a Markov chain from a Bayesian network so that the stationary distribution is meaningful?



Approximate inference using Markov chains

Bayesian network has variables
 v₁, v₂, ..., v_n



< ∃
- Bayesian network has variables
 v₁, v₂, ..., v_n
- Each assignment of values to the variables is a state



- Bayesian network has variables
 v₁, v₂, ..., v_n
- Each assignment of values to the variables is a state
- Set up a Markov chain on these states



- Bayesian network has variables
 v₁, v₂, ..., v_n
- Each assignment of values to the variables is a state
- Set up a Markov chain on these states
- Gibbs sampling random walk through state space, count visits to each state



- Bayesian network has variables
 v₁, v₂, ..., v_n
- Each assignment of values to the variables is a state
- Set up a Markov chain on these states
- Gibbs sampling random walk through state space, count visits to each state
- Stationary distribution should assign to state s the probability P(s) in the Bayesian network



- Bayesian network has variables
 v₁, v₂, ..., v_n
- Each assignment of values to the variables is a state
- Set up a Markov chain on these states
- Gibbs sampling random walk through state space, count visits to each state
- Stationary distribution should assign to state s the probability P(s) in the Bayesian network
- How to reverse engineer the transition probabilities to achieve this?

