

From last time<https://kunalmenda.com/2019/02/21/causation-and-correlation/>

How can we determine if measuring one R.V. will reveal info about another?

could. indep.

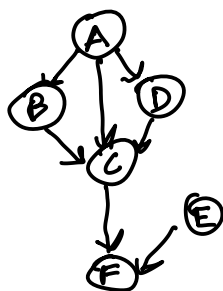
What does "Markov" mean in "MDP"?

d-separation

Def: A stochastic process is Markov w.r.t. state s iff

$$P(s_{t+1} | s_t, s_{t-1}, s_{t-2}, \dots, s_1) = P(s_{t+1} | s_t)$$

Aside: Sampling from BN



Topological Sort: if $X \rightarrow Y$ then X comes before Y

A	A
B	D
C	B
E	C
F	F
F	F

Algorithm in book

Sample A from $P(A)$

Sample B from $P(B|A)$

D from $P(D|A, B)$

⋮

Markov Chain already sorted

"simulate" a "trajectory" from Markov chain

sample s_i from $P(s_i)$

while not ended

sample s' from $P(s'|s)$

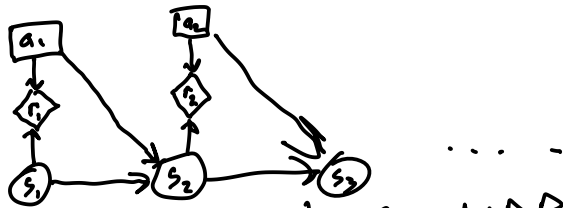
What is an MDP?
 What is a policy?



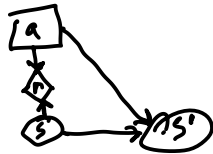
Dynamic Bayes Net



Markov Decision Process



Dynamic DN for MDP






Optimization Problem

$$\text{maximize } E \left[\sum_{t=0}^{\infty} r_t \right]$$

dynamics

Decision Network

-  decision node
-  chance
-  utility

Keeping rewards finite

1) Finite Time

$$\sum_{t=1}^T r_t$$

2) Average reward

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n r_t$$

3) Discounting γ : discount $[0,1)$

$$\sum_{t=1}^{\infty} \gamma^t r_t$$

0.9, 0.95, 0.99

if $\hat{r}_t \leq \bar{r}$

then

$$\sum_{t=0}^{\infty} \gamma^t r_t \leq \frac{\bar{r}}{1-\gamma}$$

4) Terminal / absorbing states

Infinite time but problem reaches terminal state (no reward, no leaving)
w.p. 1

"Tuple Definition"

Markov process
(S, T)

MDP defined by S-tuple

(S, A, T, R, γ)

S - state space - set of states

e.g. $\{1, 2, 3\}$, $\{\text{healthy, pre-cancer, cancer}\}$, \mathbb{R}^4

A - action space - set of actions

e.g. $\{1, 2, 3\}$, $\{\text{treat, watch, wait}\}$, \mathbb{R}^2 ,

$\{\text{working, malfunctioning}\} \times \mathbb{R}^6$

T - transition kernel

Explicit: $T(s'|s, a)$

$T: S \times A \times S \rightarrow [0,1]$ discrete
 $T(s, a, s')$

Generative: $s' = G(s, a, w)$

R - reward function

$R: S \times A \times S \rightarrow \mathbb{R}$

or $R: S \times A \rightarrow \mathbb{R}$

$R(s, a, s')$

$R(s, a)$

useful
reward for distance
traveled

$$R(s, a) = \mathbb{E}_{s' \sim T(s, a)} [R(s, a, s')]$$

γ - discount

$\gamma \in [0, 1)$

Optimization Variable

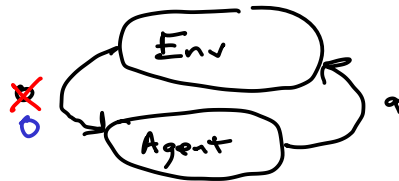
$$\underset{?}{\text{maximize}} E[\sum_t r_t]$$

Open Loop

sequence of actions

Closed Loop

choose action based on observations



In MDP observations are states

~~Policy~~
"Policy" π

$$\pi: S \rightarrow A$$

Closed Loop

$$\underset{\pi: S \rightarrow A}{\text{maximize}} E \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, s_{t+1}) \mid a_t = \pi(s_t), s_{t+1} \sim T(s_t, a_t) \right]$$

Open Loop

$$\underset{(a_1, a_2, \dots, a_T)}{\text{maximize}} E \left[\sum_{t=0}^T \gamma^t R(s_t, a_t, s_{t+1}) \mid s_{t+1} \sim T(s_t, a_t) \right]$$

Policy Evaluation

Monte Carlo Policy Evaluation

function simulate(MDP, π)

$r_{\text{sum}} = 0$ ~~sample~~ sample s from S

 for $t = 1:10,000$

$a = \pi(s)$

 sample s' from $T(s, a)$

$r_{\text{sum}} += R(s, a, s')$

$s = s'$

 end

 return r_{sum}

end

$R = 0$

for i in $1:10,000$

$R += \text{simulate}(\text{MDP}, \pi)$

end

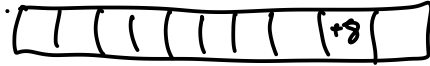
$R/N \leftarrow \text{estimate of } E[\sum_t r_t]$

Stationary Policies

$$\hookrightarrow a_t = \pi(s_t) \quad a_t = \pi_t(s_t)$$

Stationary Policies optimal

- Infinite Horizon
- Stationary Dynamics



$(S, A, T, R, \gamma, \phi)$

Example

$$T=10$$

$$\pi_1^*(\tilde{s}) = \text{right}$$

$$\pi_9^*(\tilde{s}) = \text{left}$$

