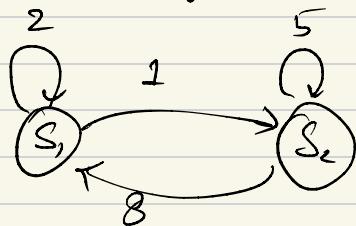


25/11/2024

## Reinforcement Learning.

$$\lambda = 0.8$$



$$D_1 = \max \{ 2 + 0.8 D_1, 1 + 0.8 D_2 \}$$

$$D_2 = \max \{ 5 + 0.8 D_2, 8 + 0.8 D_1 \}.$$

$$\text{Suppose } D_1 = 21 \text{ & } D_2 = 25.$$

$$21 = \max \{ 2 + 0.8 \times 21, 1 + 0.8 \times 25 \}$$

$$= \max \{ 2 + 16.8, 1 + 20.0 \}$$

$$= \max \{ 18.8, 21 \} = 21 //$$

$$D_2 = \max \{ 5 + 0.8 \times 25, 8 + 0.8 \times 21 \}$$

$$= \max \{ 5 + 20, 8 + 16.8 \}$$

$$= \max \{ 25, 24.8 \}$$

$$= 25 //$$

We get optimal discounted state for all states.

$D$  is function:

$$\text{state} \rightarrow [D] \xrightarrow{\quad} \text{Number(IRR)}$$

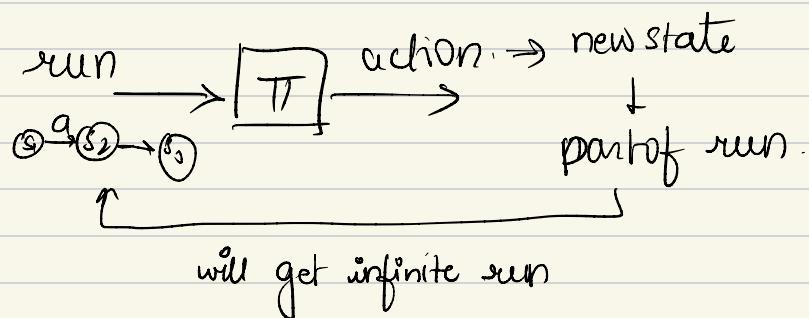
Graph  $G$ . you also have  $\text{OPT}(G) \rightarrow$  (The two  $D_1$  &  $D_2$  equations), if you find  $D_1$  &  $D_2$ , then  $D(s) = \text{optimal discounted reward from } s$ .

The policy will be memoryless and deterministic.

for every graph, there always will be a unique solution of  $\text{OPT}(G)$ .

Solution in poly time.

policy:



General policies: probability distribution over available actions.  
(but not necessary for deterministic / MDPs).

Discounted Rewards:  $r: s_0, a_1, s_1, a_2, s_2 \dots$

$$D(r) = \sum_{i=0}^{\infty} \lambda^i \cdot r(s_i, a_{i+1}).$$

$$D(s_1, \frac{b}{1}, s_2, \frac{a}{5}, s_2, \frac{b}{8}, s_1, \frac{a}{2}, s_1, \frac{b}{1}, \dots)$$

$$= 1 + \lambda^5 + 8\lambda^2 + 2\lambda^3 + \lambda^4 + \dots$$

\* Can be infinite policies.

for a state  $s$ , value or optimal discounted reward is :

$$D_*(s) = \sup_{u \in \Sigma} [D(s, u)].$$

A policy is optimum if  $D(s, u) = D_*(s)$  for every  $s \in S$ .

Proof for Correctness of Bellman Optimality Equations:

If  $D \models \text{Opt}(G)$  then  $D(s) = D_*(s)$  &  $\mu^D$  is an optimal (memoryless) policy

satisfies

**Proof:** to show:  $D(s) = \sup_{u \in \Sigma} D(s, u)$

$$s_0 \xrightarrow{a_1} s_1 \xrightarrow{a_2} s_2. \quad \text{Run}(s, u)$$

$D \models \text{OPT}(G)$ .

$$D(s_0) \geq r(s_0, a_1) + \lambda D(s_1)$$

$$D(s_1) \geq r(s_1, a_2) + \lambda D(s_2)$$

:

$$D(s_0) \geq \lim_{n \rightarrow \infty} \left( \sum_{i=0}^n \lambda^i r(s_i, a_{i+1}) + \lambda^{n+1} D(s_{n+1}) \right)$$

$$\geq \text{Disc}(s, u).$$

$$\textcircled{1} \quad D(s) \geq \text{Disc}(s, u)$$

$$\textcircled{2} \quad D(s) = \text{Disc}(s, u^D)$$

↓

∴ you are using the optimal policy  $\mu^D$ .

$$D(s_0) = \sim$$

$$D(s_1) = \sim$$

$$\therefore D(s) = \text{Disc}(s, u^D).$$

How to get optimal discounted solns (D)?

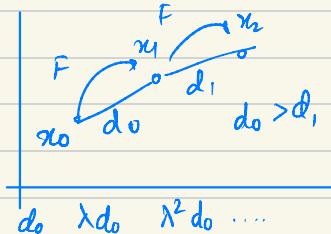
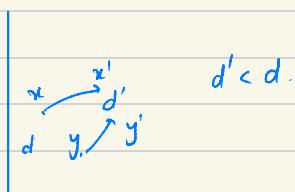
## Banach's Theorem

for every graph  $G$ , exists a unique Soln of  $\text{DPT}(G)$ .

$$D(s) = \max_{a \in A} \left\{ r(s, a) + \lambda \cdot D(T(s, a)) \right\} \quad \forall s \in S.$$

Start from  $(0,0)$   $\rightarrow (2,8) \rightarrow (7.4, 11.4) \rightarrow \dots \rightarrow (21, 25)$   
any  $(51, 3) \rightarrow \rightarrow \rightarrow$

Contraction: closeness of points. [fns that bring points closer]



$F$  is contraction if  
there exists  $\lambda \in (0,1)$   
s.t.

$$d(F(x), F(y)) \leq \lambda \cdot d(x, y)$$

Fixed point: applying the fn gives you the same point.

We can say that  $F$  is a contraction  $\Rightarrow$  fixed point is unique

$$*\text{Amazing: } x_* = \lim_{n \rightarrow \infty} x_n \quad F(x_*) = F\left(\lim_{n \rightarrow \infty} x_n\right) = \lim_{n \rightarrow \infty} F(x_n) = \lim_{n \rightarrow \infty} x_n = x_*$$

Barach's fixed point Theorem:

Every contraction mapping in a non-empty complete metric space admits a unique fixed point

all points in  
the same space      some space  
with a sense of  
distance b/w  
points -

Cauchy Sequence :

Contraction

Operator  $T$  : old value  $\rightarrow$  New value

$$v : S \rightarrow R$$

$$T(v(s)) = \max_{a \in A(s)} \{ x(s, a) + \lambda \cdot v(s) \}$$

$$s' = T(s, a).$$

To prove:

$$v_1, v_2 \quad d(T(v_1), T(v_2)) \leq \lambda \cdot d(v_1, v_2).$$

$$\text{def: } d(v_1, v_2) = \max_i |v_1(i), v_2(i)| \quad \underline{\text{max-norm.}}$$

i.  $T$  is mon

(proof given in Lec.)

$$T(v+c) = T(v) + \lambda \cdot c.$$

Theorem:  $T$  is a contraction. (To prove).

$$v_1, v_2 \quad d(T(v_1), T(v_2)) \leq \lambda \cdot d(v_1, v_2).$$

assume  $d(v_1, v_2) = c$ . &  $v_2 - c \leq v_1 \leq v_2 + c$ .

$$T(v_2 - c) \leq T(v_1) \leq T(v_2 + c).$$

$$\Rightarrow T(v_2) - \lambda c \leq T(v_1) \leq T(v_2) + \lambda c$$

$$\Rightarrow |T(v_1), T(v_2)| \leq \lambda \|v_1, v_2\|$$

$\therefore T$  is a contraction.

As a Linear Programming outline:

$$d_1 = \max(2 + 0.8d_1 + 1 + 0.8d_2)$$

$$d_2 = \max(5 + 0.8d_2 + 8 + 0.8d_1)$$

Obj: Min  $d_1 + d_2$ .

$$\text{s.t } d_1 \geq 2 + 0.8d_1$$

$$d_1 \geq 1 + 0.8d_2$$

$$d_2 \geq 5 + 0.8d_2$$

$$d_2 \geq 8 + 0.8d_1$$

$$D \geq T(D)$$

$$\Rightarrow v \geq T(v)$$

$$T(v) \geq T(T(v))$$

$$\Rightarrow v \geq T(v) \geq T^2(v) \geq \dots v_*$$

Markov Decision Process :-

Probability: finite & infinite tosses of coins:-

Random Variable: object that depends on random process.

$X$ : Sample space  $\rightarrow$  measurable  $E$ .

Expected value: probabilistic weighted average.

$$E(X) = \sum_{x \in E} x \cdot P(X=x).$$

#  $E[x_n] \rightarrow E[X]$  even if  $x_n \rightarrow x$ .

$$x_n = n \cdot 1\{U \in (0, \frac{1}{n}]\} \quad (\text{uniform dist})$$

$$x_n \rightarrow 0 \text{ but } \lim_{n \rightarrow \infty} E[x_n] = \lim_{n \rightarrow \infty} 1 = 1.$$

Monotone Convergence Theorem:

$\{X_n : n \geq 0\}$  seq of random variables with  $0 \leq X_n \leq X_{n+1}$  &  $X_n \rightarrow X$   
then  $\lim_{n \rightarrow \infty} E[X_n] = E[\lim_{n \rightarrow \infty} X_n] = E[X]$  pointwise.

Dominated Convergence Theorem :-

$\{X_n : n \geq 0\}$   $X_n \rightarrow X$  pointwise  $|X_n| \leq Y < \infty$  and  $E[Y] < \infty$   
then  $\lim_{n \rightarrow \infty} E[X_n] = E[X]$ .

## Markov Chains:

fair win to simulate a dice:

future depends only on "current".

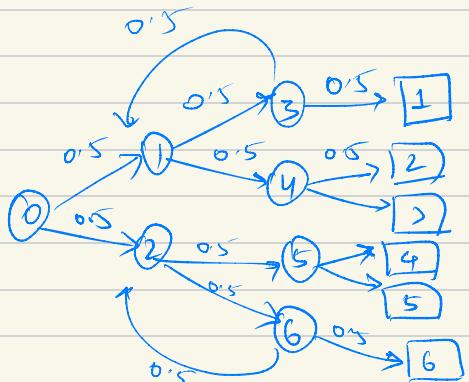
$M(S, S_0, P, r)$

$S$  is a finite set of states

$s_0 \in S$  initial state

$P : S \times S \rightarrow [0, 1]$  transition prob matrix

$r : S \rightarrow \mathbb{R}$  (reward func.)



Cylinder Sets :- (set of measurable events) [Cyl( $\omega$ )]

let finite path  $\omega = \langle s_0, s_1, s_2, \dots, s_n \rangle$  is set of all infinite paths with prefix  $\omega$ .

$$P(\text{Cyl}(\omega)) = \prod_{0 \leq i \leq k} P(s_i, s_{i+1})$$

$$P(\text{getting } 0 \text{ or } 4) = P(0, 2, 5, d^4) + P(0, 2, \underline{6}, 2, 5, d^4) + \dots$$

$$\begin{aligned} &= (0.5)^3 + (0.5)^5 + (0.5)^7 + \dots \\ &= \frac{(0.5)^3}{1 - (0.5)^2} = \frac{1}{6} // \end{aligned}$$

what does  $w$  at  
the end mean?  
(it ends not infinite) rt?

Getting the prob of a event / state reach:

$$\pi_s = \begin{cases} 1 & \text{if } s \in T \\ \sum_{t \in S} P(s, t) \pi_t & \text{otherwise.} \end{cases}$$

Consider it like the fixed point theorem for discounted values.  
 $F: [0,1]^S \rightarrow [0,1]^S$

$$F(y)(s) = \begin{cases} 1 & \text{if } s \in T \\ \sum_{t \in S} P(s, t) y(t) & \text{otherwise.} \end{cases}$$

Markov Chain:  $M(S, s_0, P, r)$

discounted reward as  $v_\lambda$   $\lambda \in [0, 1]$ .

$$v_\lambda(s) = \lim_{N \rightarrow \infty} E \left[ \sum_{t=1}^N \lambda^{t-1} r(x_t) \right] = E \left[ \sum_{t=1}^{\infty} \lambda^{t-1} r(x_t) \right]$$

Expected discounted reward:

Policy evaluation :-

~~2/4/21~~  
Sunday

## Deep RL

Learn:

- 1) value-based Methods
- 2) differences b/w Monte Carlo and Temporal Difference Learning
- 3) RL algo: Q-Learning

RL objective: Max. Rewards  $\leftarrow$  optimum policy (decision making)

2 ways:  
1) Policy-based Methods: - train agent directly to learn which action to take in a state.  
2) Value-based Methods: - train agent to learn which state is more valuable and take action that leads to it.

Value-based Methods:

Learn a function state  $\rightarrow$  expected value of being in that state.

$$\textcircled{1} \quad v_{\pi}(s) = E_{\pi}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$$

for each state,

the state value function outputs

the expected return

if the agent starts in that state.

and then follows the policy forever after.

**② Action-value function:**

$$Q_{\pi}(s, a) = E_{\pi}[G_t | S_t = s, A_t = a]$$

the action-value function outputs the expected return  
if the agent starts in that state & takes the action  
and then follows the policy forever after.

Bellman Equation: