

Markov Process (MP)

- formal description of a RL environment
  - environment is fully observable
  - current state completely characterizes the problem
  - partially observable problems can be converted to MDP
- eg: bandit problems  $\leftarrow$  MP with one state

The essential idea of a MP is :

$$P[S_{t+1} | S_t] = P[S_{t+1} | S_1, \dots, S_t]$$

$\hookrightarrow$  then  $S_t$  is a Markov state

i.e. the state is a sufficient statistic of the future

State transition probability:

$$P_{ss'} \doteq P[S_{t+1} = s' | S_t = s]$$

State transition matrix:

$$P = \begin{bmatrix} & \text{TO} \\ \text{FROM} & \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \end{bmatrix}$$

; each row sums to unity

Thus, formally, a Markov Process is a tuple  $\langle S, P \rangle$

- $S$ : finite set of states
- $P$ : state transition probability matrix

## Markov Reward Process:

- a tuple  $\langle S, P, R, \gamma \rangle$
- $R$ : immediate reward,  $R_t = E[R_{t+1} | S_t = s]$
- $\gamma$ : discount factor

Return,  $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$ ; no expectation here as it is just for one random sample from the MRP

$\gamma = 0$ : maximally short sighted

$\gamma = 1$ : maximally far sighted

## Why discount?

- mathematically convenient
- avoid infinite return in cyclic Markov processes
- uncertainty about future may not be fully represented
- may earn more from immediate rewards (say in finance)
- undiscounted MRP ( $\gamma = 1$ )

## Value function:

- long term value of a state  $s$

$$v(s) = E[G_t | S_t = s]$$

↑ as it is a stochastic process

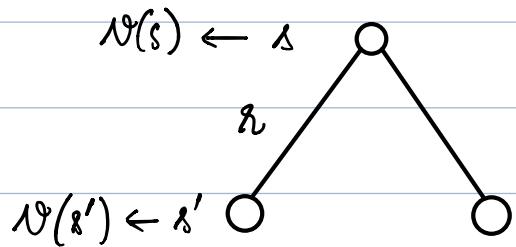
## Bellman Equation for MRP:

$$\begin{aligned}
 v(s) &= E[G_t \mid s_t = s] \\
 &= E[R_{t+1} + \gamma G_{t+1} \mid s_t = s] \quad \text{iff law of iterated expectations} \\
 &= E[R_{t+1} + \gamma v(s_{t+1}) \mid s_t = s]
 \end{aligned}$$

immediate reward      discounted value of  
 successor state

Reward for state  $s_t$  is obtained at time step  $t+1$

↳ convention to separate agent environment interface



backup diagrams  
single step lookahead search

$$\therefore v(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} v(s')$$

## Bellman Equation in Matrix form:

$$V = R + \gamma P V$$

↑      ↑      ↑  
 column vectors

$$\text{or, } (I - \gamma P)V = R$$

$$\text{or, } V = (I - \gamma P)^{-1} R$$

- $O(n^3)$  for  $n$  states
- iterative methods : DP, Monte Carlo evaluation, TD learning
  - ↳ more efficient

## Markov Decision Process:

- staple  $\langle S, A, P, R, \gamma \rangle$
- $A$ : finite set of actions
- $P$ :  $P_{s,a}^{a'} = \Pr [S_{t+1} = s' \mid S_t = s, A_t = a]$
- $R$ :  $R_s^a = E [R_{t+1} \mid S_t = s, A_t = a]$

Policy:  $\pi(a|s)$ : distribution of actions over states

$$\pi(a|s) = P[A_t = a \mid S_t = s]$$

- it defines the behaviors of an agent
- policies depend on current state only
- policies are stationary

State Value function:  $(V_\pi(s))$

$$V_\pi(s) = E_\pi [G_t \mid S_t = s]$$

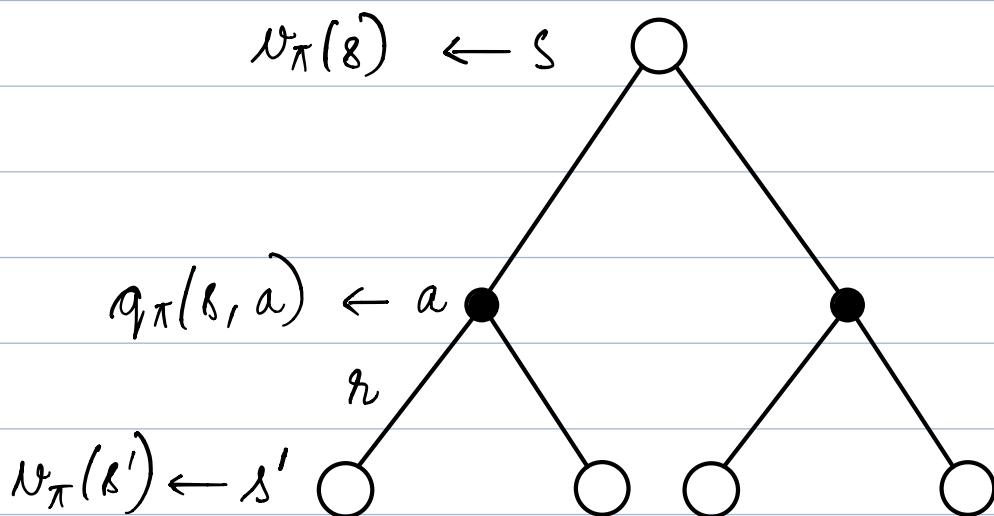
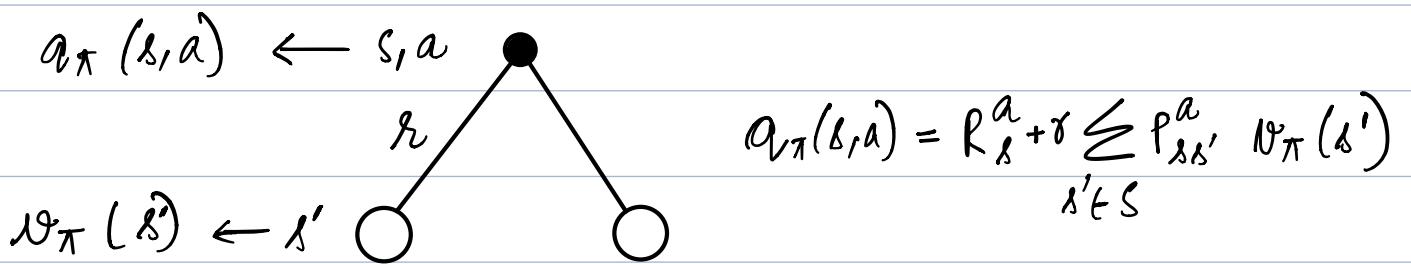
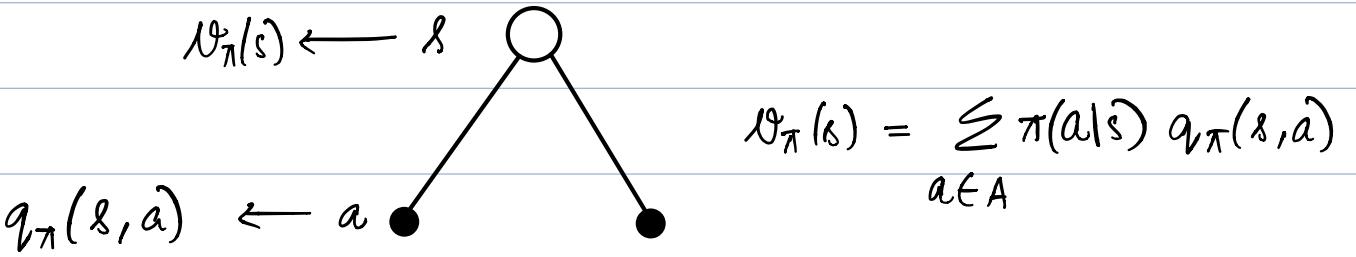
Action Value function:  $(q_\pi(s, a))$

$$q_\pi(s, a) = E_\pi [G_t \mid S_t = s, A_t = a]$$

Bellman Expectation Equation:

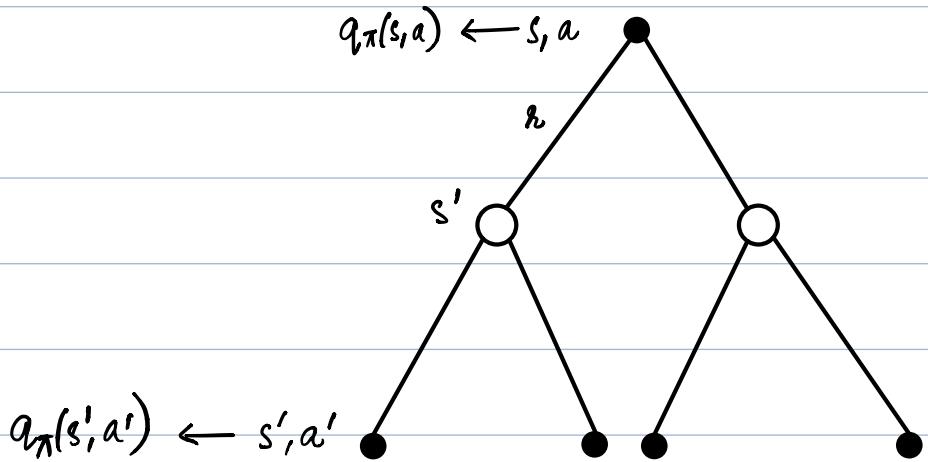
$$V_\pi(s) = E_\pi [R_{t+1} + \gamma V_\pi(S_{t+1}) \mid S_t = s]$$

$$q_\pi(s, a) = E_\pi [R_{t+1} + \gamma q_\pi(s_{t+1}, a_{t+1}) \mid S_t = s, A_t = a]$$



$$\therefore v_\pi(s) = \sum_{a \in A} \pi(a|s) q_\pi(s, a)$$

$$= \sum_{a \in A} \pi(a|s) \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_\pi(s') \right)$$



$$\therefore q_\pi(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_\pi(s')$$

$$= R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \sum_{a' \in A} \pi(a'|s) q_\pi(s', a')$$

Optimal Value function:

$$- V_\star(s) = \max_\pi V_\pi(s)$$

$$- q_\star(s, a) = \max_\pi q_\pi(s, a)$$

if this is known, then the MDP is solved

## Bellman Optimality Equation:

$V^*$ :

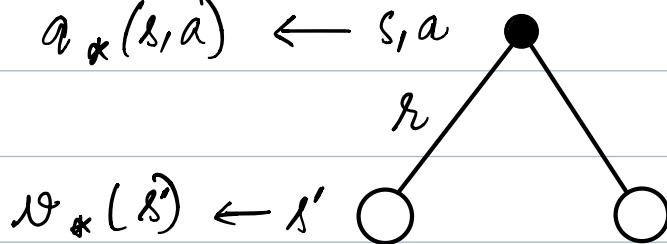
$$V^*(s) \leftarrow \delta$$

$$q^*(\delta, a) \leftarrow a$$

$$V^*(s) = \max_a q^*(s, a)$$

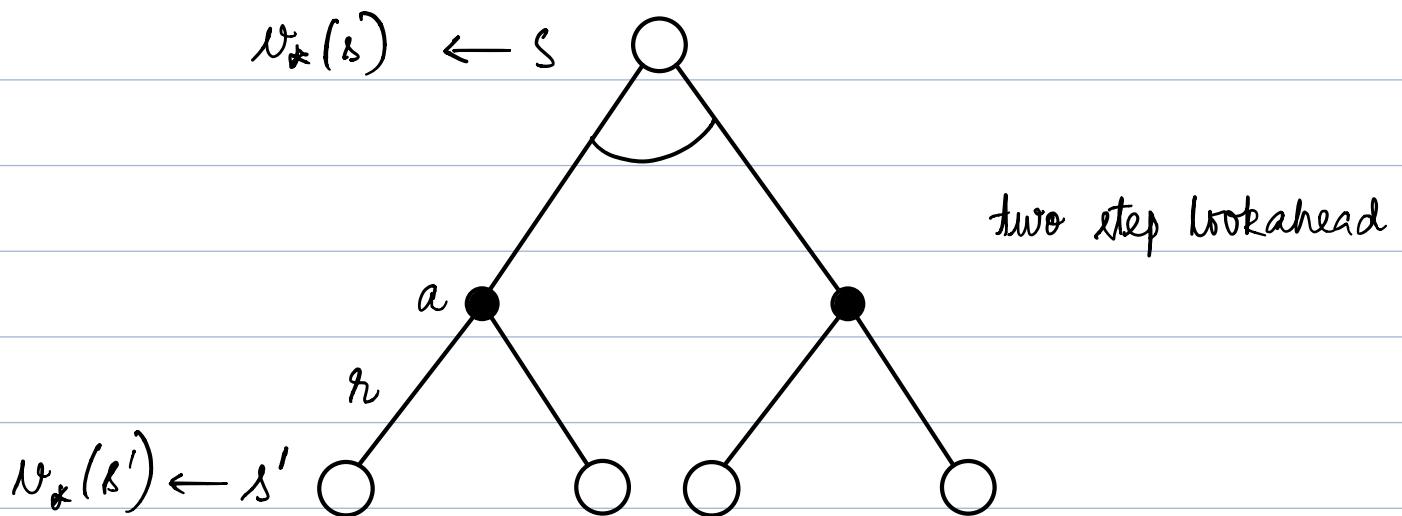
$Q^*$ :

$$q^*(s, a) \leftarrow s, a$$

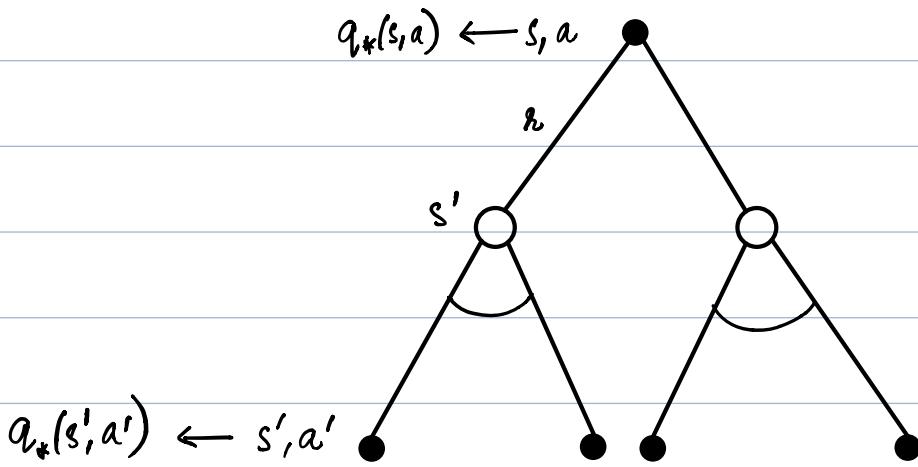


$$q^*(s, a) = r_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V^*(s')$$

$$V^*(s) \leftarrow s$$



$$\therefore V^*(s) = \max_a R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V^*(s')$$



$$q_*(s, a) = R_s^a + \gamma \sum P_{ss'}^a v_*(s')$$

$$= R_s^a + \gamma \sum P_{ss'}^a \max_{a'} q_*(s', a')$$

### Solving the Bellman Optimality Equation:

- non linear equations
- no closed form solution in general
- requires iterative solutions
  - value iteration
  - policy iteration
  - Q-learning
  - Sarsa

### Extensions to MDPs:

- infinite & continuous MDPs
- partially observable MDPs
- undiscounted, average reward MDPs