# CSPB3202 Artificial Intelligence

# MDP and RL

University of Colorado **Boulder**

Note: Some slides were adapted from the materials from Prof. Pieter Abbeel and Prof. Dan Klein at UC Berkeley, and Russell & Norvig's AI book.

# MDP: Non-deterministic Search
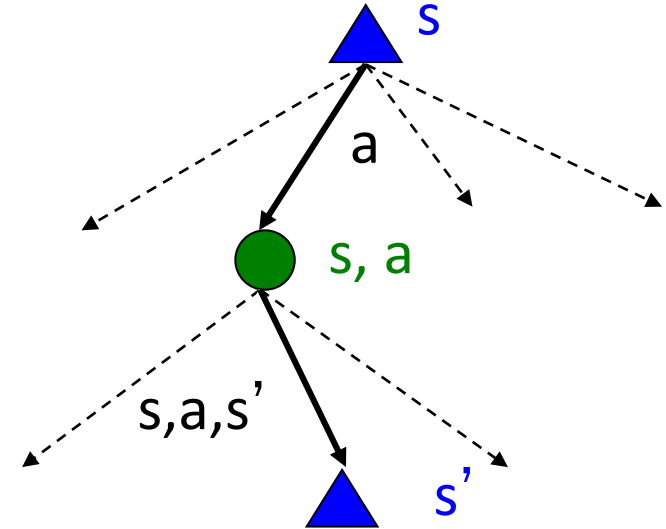
# Recap: Defining MDPs

- Markov decision processes:
  - Set of states S    $\alpha 1$: finite states
  - Start state $s_0$
  - Set of actions A    $A2$:
  - Transitions P(s'|s,a) (or T(s,a,s')) $\rightarrow$ known
  - Rewards R(s,a,s') (and discount $\gamma$)    $R(s)$    $\boxed{+1}$    living cost
    
    $\boxed{-1}$

- MDP quantities so far:
  - Policy = Choice of action for each state
  - Return (Utility) = sum of (discounted) rewards

$$G_t R = \sum_{t=0}^{\infty} \gamma^t r_t$$
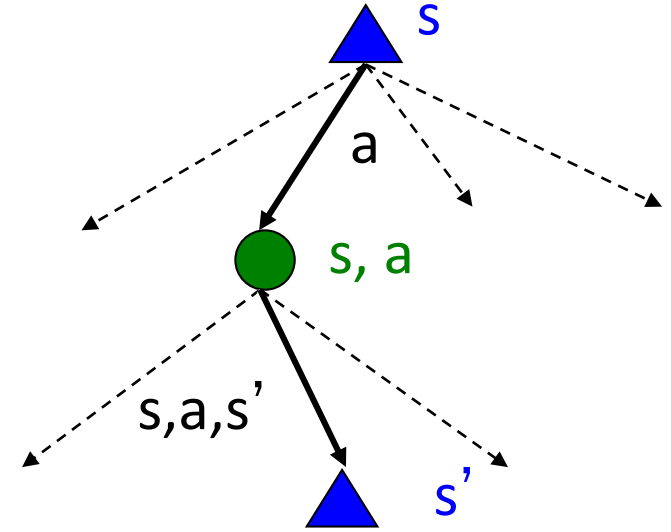
# Solving MDP

Dynamic Programming

Finite MDP — All s known (fully observable)

— s are finite

— $P(s'|s,a)$ is known, perfect

model of the world

Exactly solvable $\lceil 10^{20} \rightarrow$ Approximations

# Value of a state

- Set of states S
- Set of actions A
- Transitions P(s'|s,a) (or T(s,a,s'))
- Rewards R(s,a,s') (and discount γ)

$$G_t = \text{Return} = \sum_{k=0} \gamma^k r_t = r_0(s) + \gamma \cdot r_1 + \gamma^2 r_2 + \cdots$$

$$\underset{\text{value ft}}{\underline{V^\pi_{(s)}}} = \underset{\pi}{E}[G_t \mid S_t = s] : \text{expected return in } s$$

$$\pi : \text{policy} \Rightarrow \underline{\pi(a \mid s)} \qquad [a, \cdots \cdots]$$

action-value ft (Q-ft)

$$\underline{Q^\pi(s,a)} = \underset{\pi}{E}[G_t \mid S_t = s, A_t = a]$$

$$V^\pi_{(s)} = \sum_a \pi(a \mid s) Q^\pi(s,a)$$

$$Q^\pi(s,a) = \sum_{s',r} P(s',r \mid s,a) V^\pi(s)$$

S

a

s, a

s,a,s'

s'

# Value Functions

$$V_\pi(s) = \mathbb{E}_\pi[G_t \mid S_t = s]$$

Value ft

$$Q^\pi(s,a) = \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a]$$

Action −Value ft

$$G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$$

Return

Policy $\quad \pi(a|s)$

$$V_\pi(s) = \sum_a \pi(a|s) Q^\pi(s,a)$$

$$Q_\pi(s,a) = \sum_{s',r} P(s',r|s,a) V_\pi(s')$$



s

a

s, a

s,a,s'

s'

Goal

Sequence of action soln

$\pi(a|s_0)$

↳ s.t

$\boxed{a_0, a_1, a_2 \cdots}$

Maximize $G_t$

# Bellman Equation

$$V_\pi(s) = \underset{\pi}{E}\left[\, G_t \mid s \,\right]$$

$$G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k} = \boxed{r_t + G_{t+1}}$$

$$= r_t + \boxed{\sum_{k=1}^{\infty} \gamma^k r_{t+k}}$$

$$G_{t+1}$$

$$= \underset{\pi}{E}\left[\, r_t + \gamma G_{t+1} \mid S_t = s \,\right]$$

$$= R(s) + \gamma \underset{\pi}{E}\left[\, G_{t+1} \mid S_t = s \,\right]$$

$$s \to s'$$

$$\boxed{V_\pi(s) = R(s) + \gamma \sum_a \pi(a|s) \sum_{s'} P(s'|s,a)\, V_\pi(s')}$$

Bellman Eqn $V_\pi$

# Bellman Equation

$$\boxed{V^{\pi}(s)} = R(s) + \gamma \sum_a \pi(a|s) \sum_{s'} P(s'|s,a) \boxed{V^{\pi}(s')}$$

AIMA

U $V^{\pi}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) V^{\pi}(s')$

RL (ly SB)        $s \longrightarrow s'$

$$V^{\pi}(s) = \sum_a \pi(a|s) \sum_{s',r} P(s',r|s,a) (r + \gamma V_{\pi}(s'))$$

# Optimal Quantities

$$V^{\pi}_{(s)} = \sum_a \pi(a|s) \, q^{\pi}(s,a)$$

- **The value (utility) of a state s:**

  $V^*(s)$ = expected utility starting in s and acting optimally

  $$V^*_{(s)} = \max_{\pi} V^{\pi}_{(s)}$$

- **The value (utility) of a q-state (s,a):**

  $Q^*(s,a)$ = expected utility starting out having taken action a from state s and (thereafter) acting optimally
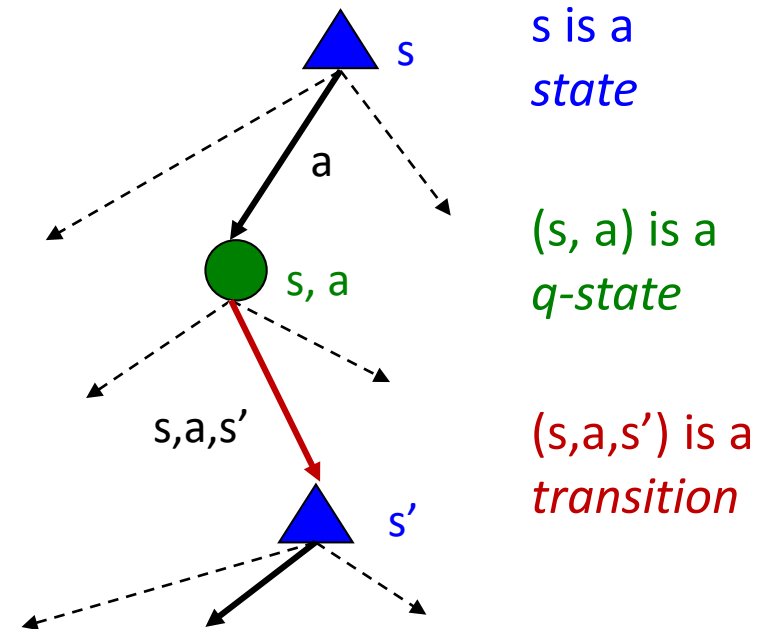
  $$q^*_{(s,a)} = \max_{\pi} q^{\pi}_{(s,a)}$$

- **The optimal policy:**

  $\pi^*(s)$ = optimal action from state s

  $$\pi^*_{(s)} = \text{argmax}_a \sum_{s'} P(s'|s,a) V_{\pi}(s')$$

s is a *state*

(s, a) is a *q-state*

(s,a,s') is a *transition*

s, a

s,a,s'

s'

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

# Bellman Optimality Equation

$$V^{\pi}(s) = \sum_a \pi(a|s) \sum_{s',r} P(s',r|s,a) \left[ r + \gamma V^{\pi}(s') \right] \qquad V^* \equiv \max_{\pi} V\pi$$

$$V^*(s) = \max_a \sum_{s',r} P(s',r|s,a) \left[ r + \gamma V^*(s') \right]$$

$$q^* \equiv \max_{\pi} q_{\pi}$$

$$V^* = \max_a q^*$$

$$q^*(s,a) = \sum_{s',r} P(s',r|s,a) \left[ r + \gamma \max_{a'} q^*(s',a') \right]$$

# Solving Bellman Equation

Dynamic Programming

1) Value Iteration

2) Policy Iteration ✓

# Solving using Policy

$$V^\pi = \sum_a \pi(a|s) \cdot$$

$$\sum_{s',r} P(s',r|s,a) \; (r + \gamma V^\pi(s'))$$

init $\begin{cases} V = 0 \\ \pi \to \text{random} \end{cases}$
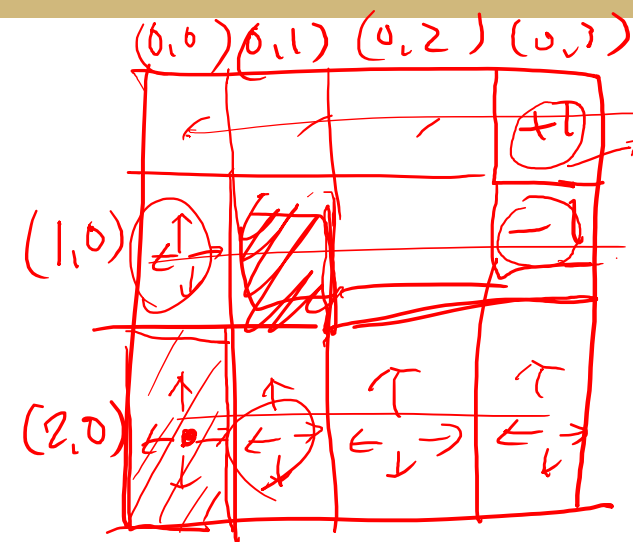
1) evaluate policy
   $V(s)$

2) Policy update

$$\max_a \sum_{s'} P(s'|s,a) V_\pi(s') \quad \sum_{s'} P(s'|s,\pi) V_\pi(s')$$

$$\pi(s) \leftarrow \operatorname*{argmax}_a \sum_{s'} P(s'|s,a) V_\pi(s')$$



Policy Iteration

$\begin{cases} \text{Eval } \pi \\ \text{update } \pi \end{cases}$

# Policy Iteration

# Value Iteration

init $V = 0 \quad \forall s$

for $i$

$\quad V(s) \leftarrow r(s) + \gamma \max_{a} \sum_{s'} P(s' | s, a) V_i(s')$

$\quad i+1$

$V_k \approx V_{k-1}$

$V_k - V_{k-1} < \Delta$

$\pi(s) = \arg\max_{a} \sum_{s'} P(s' | s, a) (r + \gamma V(s'))$

# Comparison

- Both value iteration and policy iteration compute the same thing (all optimal values)

- In value iteration:
  - Every iteration updates both the values and (implicitly) the policy
  - We don't track the policy, but taking the max over actions implicitly recomputes it

- In policy iteration:
  - We do several passes that update utilities with fixed policy (each pass is fast because we consider only one action, not all of them)
  - After the policy is evaluated, a new policy is chosen (slow like a value iteration pass)
  - The new policy will be better (or we're done)

- Both are dynamic programs for solving MDPs

# Note on Dynamic Programming

Asynchronous DP

$$\left| V_{i+1}^{(s)} - V_i^{(s)} \right|$$

$\begin{cases} \text{finite MDP} \\ \text{perfect model } (T) \end{cases}$

DP $\begin{bmatrix} V_{(s)} \rightarrow \text{Exp Return} \\ Q_{(s,a)} \quad '' \\ \pi^* \end{bmatrix}$ Bellman Equ

Policy Iter

Value Iter