

Prediction & control

Evaluate the future
- given a policy

Optimise the future
- Find the best policy

Lec-2 Markov Decision Process (MDPs)

- MDP formally describes an environment for RL, where the env. is fully observable.
i.e. the current state fully characterizes the process
- Almost all RL problems can be formalized as MDPs
- For a Markov state s and successor state s' , the state transition prob. is defined by

$$P_{ss'} = P(S_{t+1} = s' \mid S_t = s)$$

State transition matrix

to

$$P = \begin{matrix} & \text{from} & \end{matrix} \begin{bmatrix} p_{11} & \dots & p_{1n} \\ \vdots & & \vdots \\ p_{m1} & \dots & p_{mn} \end{bmatrix}$$

Each row sums to 1.

- A Markov process is a memoryless random process, i.e. a sequence of random states S_1, S_2, \dots with Markov property.

A Markov process is a tuple $\langle \mathcal{S}, P \rangle$ with

* \mathcal{S} is a (finite) set of states

* P is state-transition prob. matrix

$$P_{sr'} = P(S_{t+1} = s' \mid S_t = s)$$

- An episode is a finite sequence of Markov states when agent-env. interaction breaks naturally.

Markov Reward Process

- Markov chain with values

- MRP is a tuple $\langle \mathcal{S}, P, R, \gamma \rangle$

* \mathcal{S} set of states

* P state transition matrix

* R reward funⁿ

$$R_s = \mathbb{E}[R_{t+1} \mid S_t = s]$$

* γ is discount factor, $0 \leq \gamma \leq 1$

Return

→ Return G_t is total discounted return from time-step t

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

→ γ is present value of future rewards

- * γ trades off b/w immediate vs. long term rewards
- * Using γ provides
 - Mathematical convenient to discount rewards
 - Avoid infinite returns in cyclic Markov processes
 - Uncertainty about the future is better represented

Value function of MRP

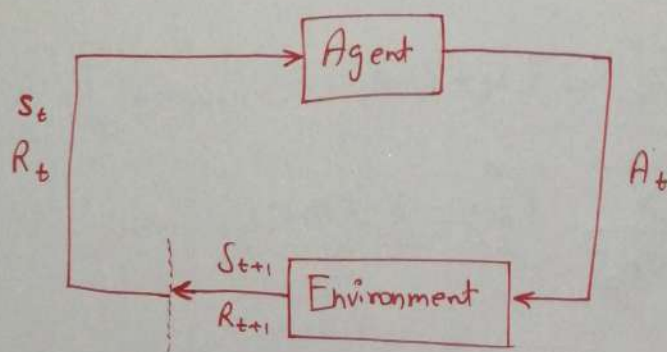
Gives the long-term value of state, State-value funⁿ of MRP

$$\begin{aligned}
 v(s) &= \mathbb{E}[G_t | S_t = s] \\
 &= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \\
 &= \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\
 &= \mathbb{E}[R_{t+1} + \gamma \mathbb{E}[G_{t+1} | S_{t+1} = s'] | S_t = s]
 \end{aligned}$$

$$v(s) = \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) | S_t = s]$$

→ from law of iterated expectation

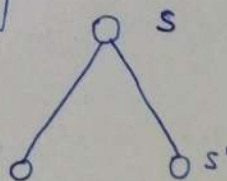
→ Bellman Equation for MRP



Agent
Environment
Interaction in
a MDP

$$v(s) = \mathbb{E}[R_{t+1} | S_t = s] + \gamma \mathbb{E}[v(S_{t+1}) | S_t = s]$$

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



→ Bellman eqⁿ for MRP

$$v = R + \gamma P v$$

$$\Rightarrow \boxed{v = (I - \gamma P)^{-1} R}$$

Markov Decision Process

An MDP is a tuple $\langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$

- * \mathcal{S} is a (finite) set of states
- * \mathcal{A} is a (finite) set of actions
- * P is state transition prob matrix

$$P_{ss'}^a = P(S_{t+1} = s' \mid S_t = s, A_t = a)$$

- * R is a reward function

$$R_s^a = E[R_{t+1} \mid S_t = s, A_t = a]$$

- * γ is a discount factor $\gamma \in [0, 1]$

Let

$$\gamma(s, a, s') = E[R_{t+1} \mid S_t = s, A_t = a, S_{t+1} = s']$$

then
$$R_s^a = \sum_{s' \in \mathcal{S}} P_{ss'}^a \gamma(s, a, s')$$

Example: Recycling Robot (Ex. 3.3 - Sutton)

- A robot with the task to collect empty cans
- State is the battery level $q = \{\text{high}, \text{low}\}$
- Actions are to either search, wait or head back to recharge

$$A(\text{high}) = \{\text{search}, \text{wait}\}$$

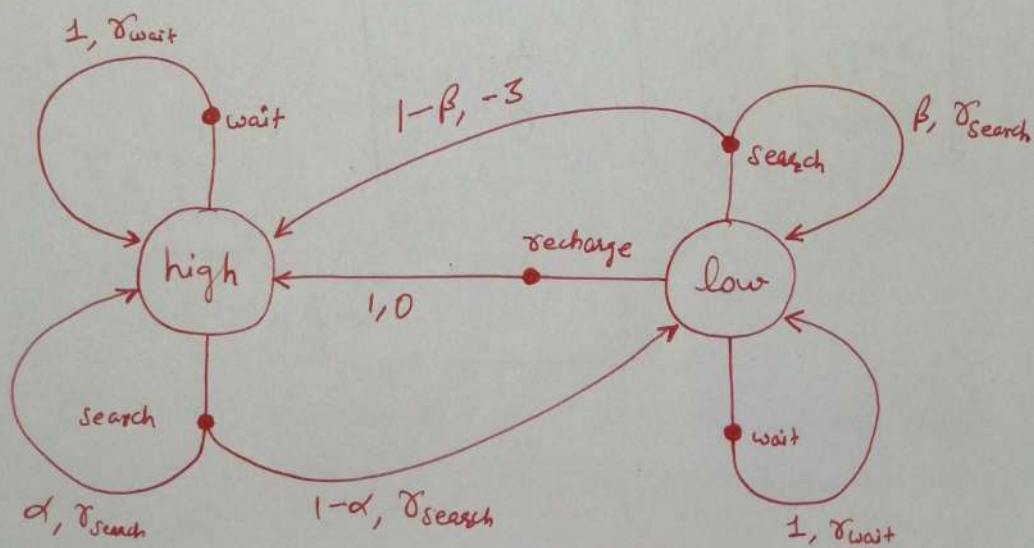
$$A(\text{low}) = \{\text{search}, \text{wait}, \text{recharge}\}$$

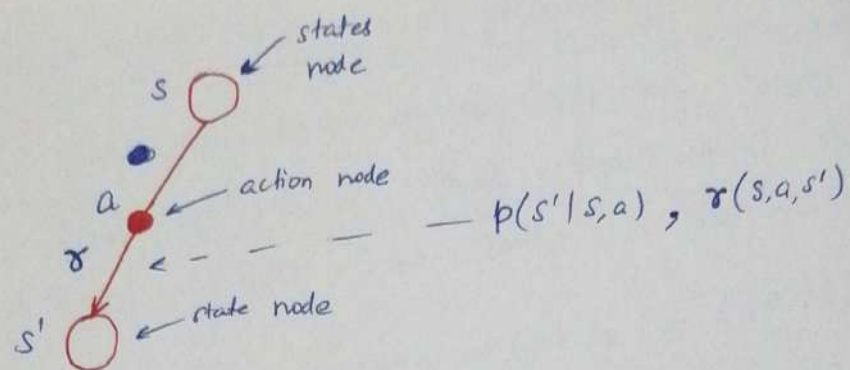
- when state is high & action is search, robot stays in high with prob. α & goes to low battery with $1-\alpha$
- similarly from low, if search action is taken, then

$$P(\text{low}) = \beta, \quad P(\text{deplete} \rightarrow \text{high}) = 1-\beta$$

- If battery depletes reward is -3 & robot is recharged to high battery
- otherwise rewards are r_{search} & r_{wait}

Transition graph for this MDP





System dynamics in tabular format

s	a	s'	$P_{s,s'}^a = p(s' s,a)$	$r(s,a,s')$
h	search	h	α	r_s
h	search	low	$1-\alpha$	r_s
h	wait	high	1	r_w
h	wait	low	0	$-$
l	search	high	$1-\beta$	-3
l	search	low	β	r_s
l	wait	high	0	$-$
l	wait	low	1	r_w
l	recharge	high	1	0
l	recharge	low	0	$-$

Policies

A policy π is a distribution over actions given states

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

A policy fully defines the behaviour of an agent.

Policies are time independent (stationary).

$$A_t \sim \pi(\cdot | s_t), \quad \forall t > 0$$

— Given an MDP $M = \langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$ and policy π

→ the state sequence S_1, S_2, \dots is a Markov process $\langle \mathcal{S}, P^\pi \rangle$

⇒ the state & reward seq. $S_1, R_1, S_2, R_2, \dots$ is an MRP

$$\langle \mathcal{S}, P^\pi, R^\pi, \gamma \rangle$$

where

$$P_{s,s'}^\pi = \sum_{a \in \mathcal{A}} \pi(a|s) P_{ss'}^a$$

$$R_s^\pi = \sum_{a \in \mathcal{A}} \pi(a|s) R_s^a$$

Value functions of an MDP

⇒ The state-value funⁿ $V_\pi(s)$ of an MDP is the expected return starting from state s , and then following policy π

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

⇒ The action-value function $q_\pi(s, a)$ is the expected return starting from s , taking action a , and then following policy π

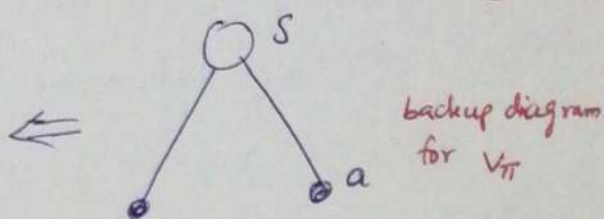
$$q_\pi(s, a) = \mathbb{E}_\pi [G_t | S_t = s, A_t = a]$$

Bellman expectation equation

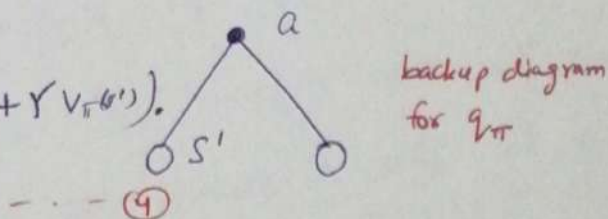
$$V_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s] \quad \text{--- (1)}$$

$$\& \quad q_{\pi}(s, a) = \mathbb{E}_{\pi} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a] \quad \text{--- (2)}$$

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) q_{\pi}(s, a) \quad \text{--- (3)}$$



$$\& \quad q_{\pi}(s, a) = \sum_{s' \in \mathcal{S}} \left(\underbrace{\delta(s, a, s')}_{P_{ss'}^a} + \gamma V_{\pi}(s') \right) \quad \text{--- (4)}$$



from (1) & (3)

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) \sum_{s' \in \mathcal{S}} P_{ss'}^a \left(\delta(s, a, s') + \gamma V_{\pi}(s') \right)$$

or

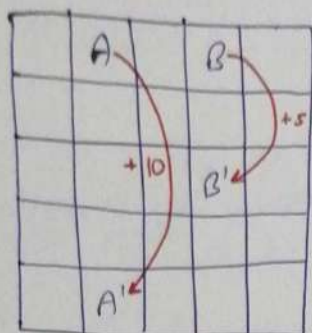
$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) \left(R_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V_{\pi}(s') \right)$$

$$\& \quad q_{\pi}(s, a) = \sum_{s' \in \mathcal{S}} P_{ss'}^a \left(\delta(s, a, s') + \gamma \sum_{a' \in A} \pi(a'|s') q_{\pi}(s', a') \right)$$

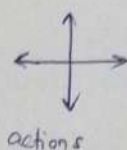
or

These are called Bellman (expectation) equations for V_π and q_π .

~~Example~~ Example: Gridworld (Sutton - Ex. 3.5)



Reward dynamics
all other ~~states~~ result in
actions
0 rewards.



		8.8		
1.5	3.0	2.3		
	0.7			

state-value funⁿ. for
equiprobable random policy
with $\gamma = 0.9$

Exercise 3.14 Let's verify that Bellman eqn holds for highlighted state.

$$\begin{aligned}
 V_\pi(s) &= 0.25 \left[0.9 \times 1.5 + 0.9 \times 8.8 + 0.9 \times 0.7 + 0.9 \times 2.3 \right] \\
 &= 0.25 \times 0.9 \times 13.3 \\
 &= 2.9925 \approx 3
 \end{aligned}$$

Exercise 3.15 Let's prove that adding a constant c to all ~~state values~~ ^{rewards} doesn't affect their relative values under any policy.

Let original rewards = r_t

& original state-value funⁿ = $V_\pi(s)$

Then new state-value funⁿ of state s

$$\bar{V}_\pi(s) = \mathbb{E} [R_{t+1} + \gamma R_{t+2} + \dots | S_t = s]$$

$$\begin{aligned}\bar{V}_\pi(s) &= \mathbb{E}[(r_{t+1} + c) + \gamma(r_{t+2} + c) + \dots | s_t = s] \\ &= \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \dots | s_t = s] \\ &\quad + \mathbb{E}[c + \gamma c + \gamma^2 c + \dots | s_t = s]\end{aligned}$$

$$\boxed{\bar{V}_\pi(s) = V_\pi(s) + \frac{c}{1-\gamma}}$$

\Rightarrow state-values of all states are incremented by constant $V_c = c/(1-\gamma)$.

Optimal Value function

The optimal state-value funⁿ is the maximum value-function over all policies:

$$\boxed{V_*(s) = \max_{\pi} V_\pi(s)}$$

similarly

$$\boxed{q_*(s,a) = \max_{\pi} q_\pi(s,a)}$$

$\Rightarrow V_*$ & q_* specify the best possible performance in the MDP, not the best policy.

\Rightarrow An MDP is "solved" when we know the optimal value functions

optimal policy

Define a partial ordering over policies

$$\pi \geq \pi' \text{ if } V_{\pi}(s) \geq V_{\pi'}(s), \forall s$$

Theorem: for any MDP, there exists an optimal policy π^* that is better than or equal to all other policies. i.e.

$$\pi_* \geq \pi, \forall \pi$$

\Rightarrow All optimal policies achieve the same optimal value funⁿ & same optimal action-value funⁿ

$$V_{\pi_*}(s) = V_*(s)$$

$$q_{\pi_*}(s) = q_*(s)$$

An optimal policy can be found by maximizing over $q_*(s, a)$.

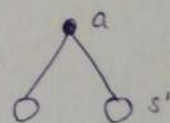
$$\pi_*(a|s) = \begin{cases} 1, & \text{if } a = \operatorname{argmax}_{a \in A} q_*(s, a) \\ 0, & \text{o/w} \end{cases}$$

\Rightarrow There is always a deterministic optimal policy for any MDP.

Bellman optimality equations

$$V_*(s) = \max_a q_*(s, a)$$

$$q_*(s, a) = \sum_{s' \in \mathcal{S}} P_{ss'}^a (r(s, a, s') + \gamma V_*(s'))$$



Now,

$$V_*(s) = \max_a \sum_{s' \in \mathcal{S}} P_{ss'}^a (r(s,a,s') + \gamma V_*(s'))$$

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

$$Q_*(s,a) = \sum_{s' \in \mathcal{S}} P_{ss'}^a (r(s,a,s') + \gamma \max_{a'} Q_*(s',a'))$$

$$Q_*(s,a) = R_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a \max_{a'} Q_*(s',a')$$

- Bellman optimality eqn is non-linear and it doesn't have any closed form solution in general.
- There are iterative methods to solve it
 - Value iteration
 - Policy iteration
 - Q-learning
 - Sarsa

Example 3.9 (Sutton): Bellman optimality eqⁿ for the recycling robot

$$V_*(h) = \max \begin{cases} P(h|h,s) [\gamma(h,s,h) + \gamma V_*(h)] + P(l|h,s) [\gamma(h,s,l) + \gamma V_*(l)] & \text{(search)} \\ P(h|h,w) [\gamma(h,w,h) + \gamma V_*(h)] + P(l|h,w) [\gamma(h,w,l) + \gamma V_*(l)] & \text{(wait)} \end{cases}$$

$$= \max \begin{cases} \alpha (\gamma_s + \gamma V_*(s)) + (1-\alpha) (\gamma_s + \gamma V_*(l)) \\ \gamma_w + \gamma V_*(h) \end{cases}$$

$$= \max \begin{cases} \gamma_s + \gamma (\alpha V_*(s) + (1-\alpha) V_*(l)) \\ \gamma_w + \gamma V_*(h) \end{cases}$$

similarly

$$V_*(l) = \max \begin{cases} \beta \gamma_s - \beta(1-\beta) + \gamma [(1-\beta) V_*(h) + \beta V_*(l)] \\ \gamma_w + \gamma V_*(l) \\ \gamma V_*(h) \end{cases}$$

for any choice of $\gamma_s, \gamma_w, \alpha, \beta \in \gamma$ with $0 \leq \alpha, \beta \leq 1$, $0 \leq \gamma \leq 1$, there is exactly one pair of numbers $V_*(h) \triangleq V_*(l)$ that simultaneously satisfy these two nonlinear eqⁿs.