

Ch-3, Finite Marker Decision Processes

→ Marker Decision Process:

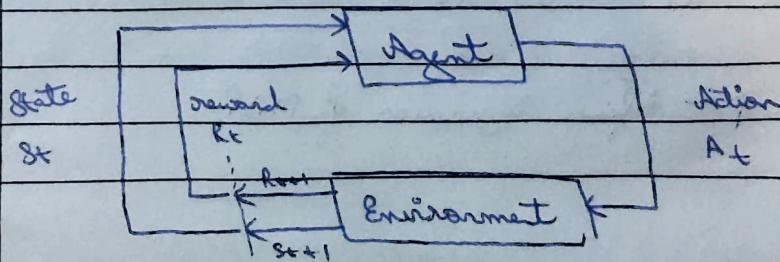
- Classical formalization of sequential decision making, where actions not only influence immediate rewards but also subsequent situations or state and future rewards.
- i.e. MDP involves delayed reward and need to tradeoff for immediate and delayed rewards.
- we estimate value $q_*(s, a)$ of each action a in each state s , or we estimate the value $v_*(s)$ of each give optimal action selection.
- Idealized form of RL problems for which precise theoretical statements can be made.

The Agent - Environment Interface :-

→ Agent - Learner and decision maker

Environment - The thing agent interacts with, comprising everything outside the agent.

→ They interact continuously, the agent selecting actions and the env responding with new situations and rewards.



Agent - Env interaction
in a Marker Decision process

- The interaction b/w the agent and the env occurs at discrete time steps ($t=0, 1, 2, 3 \dots$)
 - At each time step t , the agent receives some representation of the env's state, $s_t \in S$ and on that basis selects an action, $a_t \in A(s)$.
 - One time step later, in part as a consequence of its action, the agent receives a numerical reward $r_{t+1} \in R$ and finds itself in a new state, s_{t+1} .
 - This iteration generates a sequence of states, action and rewards :
- $$s_0, a_0, r_1, s_1, a_1, r_2, s_2, a_2, r_3 \dots$$
- In finite MDP, sets of actions, states and rewards (S, A, R) have a finite number of elements.

r_t, s_t - random variables with well defined discrete prob distributions dependent only on the preceding action and state

$$P(s', r|s, a) \doteq P_r \{ s_t = s', r_t = r | s_{t-1} = s, a_{t-1} = a \}$$

↑
for particular values of $R, s \in S$ and $a \in A$, there is a prob of those values occurring at tie t , given particular value of the preceding state and action.

for all $s, s' \in S, a \in A$ and $a \in A(s)$

• future P - defines dynamics of MDP

- The dot over the equal sign in the eqⁿ reminds us that it is a definition rather than a fact that follows from previous definition.
 - Dynamic function $p: S \times R \times S \times A \rightarrow [0, 1]$ is a ordinary deterministic function of 4 arguments.
- $\rightarrow \sum_{s' \in S} \sum_{a \in A} p(s', r|s, a) = 1$ for all $s \in S, a \in A(1)$
- p - Probability distribution for each choice of S and a
- \rightarrow In Markov decision process, the prob given by p completely characterize the env's dynamics

That is, prob of each possible value for s' at $t+1$ depends only on the state immediately preceding state at action, s_t at $t-1$, and given them, not at all on earlier states or actions.

This can be best viewed a restriction not on the decision process, but on the state.

- * State must include if all aspects of the past agent-env interaction that makes a difference for the future. Markov property.
- \rightarrow From the 4 arg dynamic function p , one can compute anything else one might want to know about the env, such as the state-transition prob (which we denote, with a slight abuse of notation, as 3 arg function $p: S \times S \times A \rightarrow [0, 1]$)

$$p(s' | s, a) = p_{\pi}(s_t = s' | S_{t-1} = s, A_{t-1} = a)$$

$$= \sum_{a \in R} p(s'_t, a | s, a)$$

Compute reward for state-action pair as a 2 argument function $\pi: S \times A \rightarrow R$:

$$\pi(s, a) \doteq E[R_t | S_{t-1} = s, A_{t-1} = a]$$

$$= \sum_{a \in R} \sum_{s' \in S} p(s'_t, a | s, a)$$

The expected reward for state-action-next state triple as a 3 arg function $\pi: S \times A \times S \rightarrow R$

$$\pi(s, a, s') = E[R_t | S_{t-1} = s, A_{t-1} = a, S_t = s']$$

$$= \sum_{a \in R} \frac{p(s'_t, a | s, a)}{p(s' | s, a)}$$

→ MDP framework - abstract & flexible so can be applied to many different problems in many different ways.

• The steps, action, states can be defined in various way to suit different tasks.

Note: Actions can be any decisions we want to learn how to make.

States can be anything we know that might be useful in making these

- The boundary between the agent and the environment is determined by what agent can control, anything beyond the agent's control is part of env.
- Reward computation to be external of an agent because it defines the task facing the agent and thus must be beyond its ability to change arbitrarily.
- MDP framework reduces the problem of learning goal-directed behavior to 3 signals:
 - the actions chosen by agent
 - states that inform these choices
 - rewards that define agent's goal

Ex

Pick and Place Robot:

- It has to control motor of a robot arm
- Learning agent will have to control the motors directly and have low latency info about current pos and velocity of mechanical linkages.
- Action - voltage applied to each motor at each joint
- State - latest reading of joint angles and velocities
- Reward - +1 for each obj picked up successfully at place
- To encourage smooth movement, -ve reward can be given as a function of the moment-to-moment "jerkiness"

Note

The action depends on the state it is in.

The transition prob label by the arrow leaving a state node always sum to 1

Real Pg 2-3 (Pengui Robots) Pg - 52-53

Goals and Rewards :-

- The purpose or goal of the agent is formalized in terms of a special signal called reward, passing from environment to agent.
- At each time step the reward is a single number R_t .
- Agent's job: Maximizing total amt of reward it receives over time, focusing not just on immediate rewards but on cumulative rewards in long run.
- Reward Hypothesis: Suggests that all goals and purposes can be effectively represented as the maximization of the expected value of the cumulative sum of received rewards.
- The reward signal is your way of communicating to the robot what you want it to achieve, not how you want it achieved.
- The reward signals can be tailored to represent different goals as objectives in RL framework.

Returns and Episodes :-

- Agent's goal: Maximizing reward it receives in long run.
This is formally defined in terms of the expected return, where the return (G_t) is a function of the sequence of rewards received after the step t .

Return = Sum of rewards

$$\text{Return} \rightarrow R_T = R_{t+1} + R_{t+2} + \dots + R_T$$

where T - final time step

- This formulation is suitable for episodic tasks, where the agent-environment interaction naturally breaks into subsequences or episodes, such as plays of game or any sort of repeated interaction.
- Each episode ends in a terminal state, followed by a reset to a standard starting state.
- Episodes are considered to end in the same terminal state with different rewards for different actions.
Tasks with episodes of this kind - episodic tasks.
- Set of all non-terminal states - S
Set of all states plus terminal state - S^+
- T - random variable that can vary from episode to episode
- However, for continuing tasks, where the interaction does not break into identifiable episodes all goes on continually without limit, the return formulation is problematic because time step $T = \infty$ at the return itself could be ∞ .
- To address this, the concept of discounting is introduced.
- Discounting involves weighty rewards received in the future less than immediate rewards, which ensures that the sum of

The discounted ~~gains~~ rewards remains finite even in
Continuing Tasks

The agent then seeks to maximize the expected discounted return:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

where γ - a discount rate ($0 \leq \gamma < 1$)

- γ determines the present value of future rewards: a reward received k time steps in future is worth only γ^k times ~~less~~ what it would worth if it were received immediately.
- If $\gamma < 1$, the ∞ sum has a finite value as long as reward sequence $\{r_k\}$ is bounded
- $\gamma = 0$, the agent - myopic - only concerned with immediate rewards - objective to learn how to choose so as to only maximize R_{t+1} .
- γ - reflects tradeoff b/w immediate future rewards, with a lower γ valuing immediate rewards more heavily.
This formulation allows the agent to make decisions that balance short term gains with long term benefits, even in tasks that continue indefinitely.

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

$$G_t = R_{t+1} + \gamma G_{t+1}$$

Finite Return is a sum of an ∞ no of terms, it is still finite if reward is non-zero and constant if $\gamma < 1$
 For e.g., if reward is a const +1, then return is

$$G_t = \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}$$

Unified Notation for Episodic and Continuing Tasks :-

→ Tasks can be categorized into episodic tasks, where the agent - env interaction breaks down into separate episodes and continuous tasks where interaction does not break into distinct episodes.

Episodic tasks are simpler because each action affects only a finite number of subsequent rewards within the episode

→ For episodic tasks, the steps are numbered starting anew from 0 for each episode

∴ $s_{t,i}$ - state representation at time t of episode i

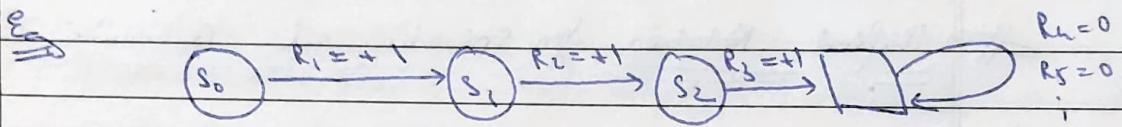
Similarly $a_{t,i}$, $r_{t,i}$, T_i , $\pi_{t,i}$ etc.

However in practice, the explicit reference to the episode number is often omitted when discussing episodic tasks, as the focus is usually on a single episode or statements that apply to all episodes.

That is we write s_t to refer to $s_{t,i}$

→ To unify the notation for episodic and continuing tasks, return is defined as a sum over an infinite number of steps, even for episodic tasks.

This is achieved by considering episode termination as entering a special absorbing state that transitions only to itself and gives only rewards of 0.



- Square box - represents the special absorbing state corresponding to the end of an episode
- Starting from s_0 , we get reward seq $+1, +1, +1, 0, 0, 0, \dots$
- Running these, we get same return whether we run over the first T rewards (here $T=3$) or over full infinite seq.
- This remains true even if we introduce discounting

∴ The return can be defined using the same eqⁿ for both types of tasks, with the understanding that the sum remains finite for episodic tasks due to absorbing state.

The return for both episodic and continuing tasks can be written as:

$$G_t = \sum_{k=t+1}^T \gamma^{k-t-1} R_k,$$

Where T - final time step for episodic tasks (which may be ∞ for continuing tasks)

γ - discount factor

$T=\infty$ or $\gamma=1$ (but not both)

- This convention simplifies the notation and highlights the parallels b/w episodic and continuing tasks in RL

Policies and Value Functions :-

- RL often involve estimating value functions, which are functions of states or state-action pairs that estimate how good it is for agent to be in a given state or perform a given action in a given state.
- "Goodness" - defined in term of the expected future or rewards or the the expected return
- value functions are defined wst specific ways of acting - policies
- Policy - mapping from states to possibilities of selecting each possible action
- If the agent is following policy π at the tth, the $\pi(a|s)$ is the prob that the action A is taken if the state St is S.
- RL's methods specify how the agent's policy is changed as a result of its experience
- value function of a state S under a policy π , denoted by $V\pi(s)$, is the expected return when starting in S and following π thereafter.

$$v_{\pi}(s) \doteq E_{\pi} \left[h_t \mid S_t = s \right]$$

$$= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right] \text{ for all } s$$

where E_{π} - expected return of a random variable given that the agent follows policy π and t is any time step

v_{π} - state value function for policy π

Note Value of the terminal state = 0

→ Similarly, value of taking action a in state s under a policy π , denoted by $q_{\pi}(s, a)$ as the expected return starting from s , taking action a and thereafter following policy π :

$$q_{\pi}(s, a) = E_{\pi} \left[h_t \mid S_t = s, A_t = a \right]$$

$$= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

q_{π} - action value function for policy π

→ Value functions q_{π} and v_{π} can be estimated from experience

For e.g.: if an agent follows policy π and maintains a avg for each state encountered of the actual returns that have followed that state, then the avg will converge to the state's value, $v_{\pi}(s)$, as the no. of times that state is encountered approaches ∞ .

Similarly keeping separate avg for each action taken in each

state, then these averages will similarly converge to the action values $q_{\pi}(s, a)$

These estimation methods - Monte Carlo methods - they involve averaging over many random samples of actual returns

* → Value function satisfy recursive relationships known as Bellman eqⁿ.

For any policy π and any state s , the Bellman eqⁿ for $V_{\pi}(s)$ is:

$$\begin{aligned} V_{\pi}(s) &= E \{ h_t \mid S_t = s \} \\ &= E_{\pi} \{ R_{t+1} + \gamma h_{t+1} \mid S_t = s \} \\ &= \sum_{a \in A} \pi(a \mid s) \sum_{s' \in S} \sum_{r \in R} p(s', r \mid s, a) [r + \gamma E_{\pi} \{ h_{t+1} \mid S_{t+1} = s' \}] \end{aligned}$$

$$V_{\pi}(s) = \sum_{a \in A} \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma V_{\pi}(s')], \text{ for all } s \in S$$

where action a taken from set $A(s)$

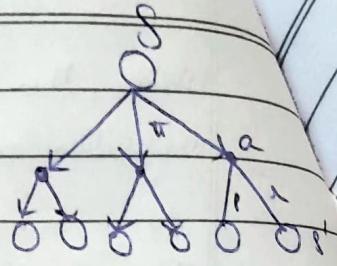
next state s' from set S

rewards r from set R

* Bellman eqⁿ expresses a relationship between the value of a state and the values of its successor state.

• Each open circle - State

Solid circle - State-action pair



Starting from state S , the robot node at top, the agent could take any of some set of action - π (from diag) based on policy π . From each of these, the env could respond with one of several next states s' along with a reward r , depending on dynamics given by P . Bellman eq¹ avg over all possibilities weighted by its prob of occuring.

It states that the value of the start state must equal the (discounted) value of the expected next state, plus the reward along the way.
(approx)

v_π - unique sol¹ to its Bellman eq¹

See the eq from
TB Pg-60-61

Optimal Policies and Optimal Value Functions:

→ Solving a RL task means, finding a policy that achieves a lot of reward over the long run.

→ for finite MDP, we define an optimal policy (π^*) that is better than or equal to all other policies.

This means that for all states s , the expected return is greater than or equal to that of π' for all state-optimal policies.

In other words, $\pi \geq \pi'$ if and only if $v_\pi(s) \geq v_{\pi'}(s)$ for all policies.

- Although there may be more than one, we denote all optimal policies by π^* .

They share same state value function called optimal state value function (v^*):

$$v^*(s) = \max_{\pi} v_{\pi}(s) \quad \text{for all } s \in S$$

- Optimal policy also ~~shares~~ ^{shares} same optimal action value function (q^*):

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

for all $s \in S$ at a GAI(s).

- For the state action pair (s, a) this function gives the expected return for taking action a in state s and thereafter following π^* an optimal policy.
- This, q^* is true of v^* :

$$q^*(s, a) = E \left\{ R_{t+1} + \gamma v^*(s_{t+1}) \mid S_t = s, A_t = a \right\}$$

- Because v^* - value function for a policy, it must satisfy the self-consistency condition given by the Bellman eq^t for state values

The Bellman optimality eq^t for v^* expresses the fact that value of a state under an optimal policy must equal the expected return for the best actions from that state:

$$v_*(s) = \max_{a \in A(s)} q_{**}(s, a)$$

$$= \max_a E_{\pi^*} \left[r_t + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a \right]$$

$$= \max_a E_{\pi^*} \left[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a \right]$$

$$= \max_a E \left[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a \right]$$

$$\rightarrow v_*(s) = \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_*(s')]$$

Bellman Optimality eqⁿ

\rightarrow Similarly, Bellman optimality eqⁿ for q_* is

$$q_*(s, a) = E \left[r_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right]$$

$$= \sum_{s', r} p(s', r \mid s, a) [r + \gamma \max_{a'} q_*(s', a')]$$

\rightarrow For finite MDP, the Bellman optimality eqⁿ for v_* has a unique solⁿ independent of policy.

The Bellman optimality eqⁿ - actually a sys of eqⁿ, one for each state, so if there are n states, then there are n equations in n unknowns.

If the dynamics p of the env are known, then in principle one can solve the sys of eqⁿ for v_* using any of a variety of methods for solving sys of non linear eqⁿ.

→ Once we have the optimal value function v^* or optimal action value function q^* , we can easily determine an optimal policy π^* by selecting actions that maximize these functions.

For any state s , an action a is optimal if it achieves the max in the Bellman optimality eq^{*}.

→ Beauty of v^* (optimal state-value) lies in its ability to simplify the decision making process. v^* is used to evaluate the short term consequences of actions, you can make decisions based on immediate info while ensuring long term optimality. This is because v^* incorporates the reward consequences of all possible future behavior.

→ Having q^* (optimal action-value) simplifies the process of choosing optimal actions even further than using the v^* .

With q^* , the agent does not need to perform a one-step ahead search to determine the best action. Instead, for any given state s , the agent can directly select an action that maximizes $q^*(s, a)$.

By using q^* , the agent trades off the simplicity of only needing to represent a function of states (as in v^*) for the convenience of being able to select optimal actions directly, without further computation. This makes q^* a powerful tool for decision making in RL.

→ In practice finding the exact solⁿ to these eq^{*} can be challenging due to computational complexity and the need for complete knowledge of the env's dynamics. ∴ RL algo

often aim to find appropriate "sol" to these eq" using
techniques like dynamic programming, Monte Carlo methods,
and temporal difference learning.

Optimality and Approximation:-

- In general learning an optimal task is a rare occurrence due to extreme computational cost associated with generating optimal policies for complex tasks.
- Even with a completely accurate model of the env's dynamics, it is usually not possible to simply compute an optimal policy by solving the bellman optimality eq.
For eg: board games like chess are a tiny fraction of human experience, yet large, custom designed computers still can't compute the optimal moves. A critical aspect of the problem facing the agent is always the computational power available to it, in particular, the act of computation can perform in a single time step.
- Computational power and memory are critical constraints for all agents.
A large amount of memory is often required to build up approximations of value functions, policies and models.
~~simplifying kept the model~~
In case with small, finite state sets - tabular method stores values in arrays or tables for each state-action pair.
- However, in many practical situations, the no of states for exceeds what can be feasibly stored in a table, necessitating the use of more compact parameterized function representations to approximate the functions.

→ RL interestingly involves settling for approx due to these constraints.

However, it also presents us with some unique opportunities for achieving useful approximation.

For Eg: an agt might encounter ~~many~~ states with such low prob that selecting suboptimal optimal actions for them has little impact on the act of reward the agt receives

This allows RL to focus more on effort on learning good decisions for frequently encountered ~~other~~ states while expending less effort on infrequently encountered states.

This selective focus on important states is a key prop that distinguishes RL from other methods for approximately solving MDP.

Read Summary
Pg. 68