

COS 435/ECE 433 Week 2 Precept Notes

February 14, 2024

1 Review

In this class, we are primarily concerned with methods for solving *Markov Decision Process* (MDP). As such, it is important to keep in mind the definition(s) of MDPs:

Definition 1.1 (Markov Decision Process, Finite-Horizon Setting). A *finite-horizon MDP* is a tuple $(\mathcal{S}, \mathcal{A}, p, r, T)$, where

- \mathcal{S} is the set of *states*,
- \mathcal{A} is the set of *actions*,
- p is the *transition dynamics*, such that $p(s'|s, a)$ is the distribution over the next state s' , given the agent took action a from state s .
- r is the *reward function*, such that $r(s, a)$ is the reward from taking action a from state s .
- T is the *rollout horizon*, i.e. the number of actions an agent takes in every rollout. \diamond

Definition 1.2 (Markov Decision Process, Discounted Infinite-Horizon Setting). A *finite-horizon MDP* is a tuple $(\mathcal{S}, \mathcal{A}, p, r, \gamma)$, where $\mathcal{S}, \mathcal{A}, p, r$ follow the same definitions as in the finite-horizon setting, and γ is the *discount factor*. \diamond

These definitions capture the notion of a particular task that an agent is trying to solve, as well as a way of measuring how well an agent is solving a task at a particular timestep (as expressed by the reward function). To quantify how well an agent is solving a task overall, we define the notion of a (*discounted*) *return*:

Definition 1.3 (Return). Let π be a policy (a way of acting in the MDP), defined as a distribution over actions conditioned on a state. Then, the return of the policy π from a state s , denoted $V^\pi(s)$, is defined as

$$V^\pi(s) := \mathbb{E} \left[\sum_{t=0}^{T-1} r(s_t, a_t) \right]$$

in the finite-horizon setting, or

$$V^\pi(s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$

in the infinite-horizon setting. \diamond

2 Discussion Problems

1. Decide whether the following scenarios are better modeled as a contextual bandit problem or a general RL problem:
 - (a) Recommendation systems (e.g. Netflix)
 - (b) A hypothetical system for recommending medical interventions over the course of a treatment plan.

[Recommender systems] Recommender systems are often modeled as contextual bandit problems, especially when the goal is to maximize immediate engagement or satisfaction. In a contextual bandit setting, the system selects actions (e.g., recommendations) based on the current context (e.g., user profile, browsing history) without explicitly considering the long-term impact of these actions. This model is suitable for scenarios where the immediate reward is the primary focus, and the effect of the current action on future states is either negligible or not considered.

[Recommending medical interventions over a treatment plan] In this case, the problem is inherently sequential and involves significant considerations for the future state of the patient. Each intervention has consequences that affect the patient's health, which influences the appropriateness of future interventions. This is a classic scenario for general RL, where each action (medical intervention) affects the environment (patient's health state) in a way that must be considered for future decisions. The objective is not just to maximize the immediate outcome of an intervention but to optimize the patient's health over the entire treatment plan. This requires a model that can account for the temporal dependencies between actions and their long-term effects on the state, making general RL the suitable approach.

2. Formulate finding the shortest path between two nodes on a given graph $G = (V, E)$ as an RL problem. Give three equivalent ways: (1) letting the state be the current node, (2) having the state maintain a stack of visited states (3) similar state space as (2), but computing the reward using the stack only at termination (sparse reward).

[(1) State = Current Node]

- States: Each state represents the current node the agent is located at within the graph G .
- Actions: Actions correspond to moving from the current node to an adjacent node.
- Transitions: Determined by the graph structure; moving from one node to another changes the state based on the edge E between them.
- Reward: A negative reward (e.g., -1) for each move encourages the agent to find the shortest path (recall that in RL we aim to maximize total reward). A large positive reward could be given upon reaching the destination node, incentivizing completion.
- Horizon: We can set the horizon to $|V|$, the number of nodes in the graph, since no shortest path will visit a node twice.

[(2) State = Stack of Visited States]

- States: State maintains the history of visited states so far.
- Actions: Actions correspond to moving from the current node to an adjacent node.
- Transitions: Determined by the graph structure; moving from one node to another adds a new node to the state stack.

- **Reward:** A negative reward (e.g., -1) for each move encourages the agent to find the shortest path (recall that in RL we aim to maximize total reward). A large positive reward could be given upon reaching the destination node, incentivizing completion.
- **Horizon:** We can set the horizon to $|V|$, the number of nodes in the graph, since no shortest path will visit a node twice.

[(3) State = Stack of Visited States, Sparse Rewards]

- **States:** State maintains the history of visited states so far.
- **Actions:** Actions correspond to moving from the current node to an adjacent node.
- **Transitions:** Determined by the graph structure; moving from one node to another adds a new node to the state stack.
- **Reward:** Reward is only given once the stack length reaches $|V|$ or when the agent reaches the goal state, whichever one comes first. The reward is computed using the stack (and computed in a way to be equivalent with the earlier two formulations).
- **Horizon:** We can set the horizon to $|V|$, the number of nodes in the graph, since no shortest path will visit a node twice.

Remark 2.1. Observe that there are multiple ways to formulate the same MDP. However, as we will learn throughout this course, the way we formulate the MDP can have a huge impact on how easy/difficult it is to find an optimal policy in an MDP. In the example above, note that the first formulation is straightforward and models the problem directly in terms of the graph's structure. The continuous reward signal (negative reward at each timestep) helps guide the agent by providing immediate feedback on its actions, making the learning process potentially easier. On the other hand, the latter two formulations consider a much larger state space (with $n!$ possible states worst case, depending on the graph), and the final formulation has extremely sparse reward. \lrcorner

3. (Markov chain, return computation practice) Consider an MDP where the state space is structured as a binary tree of depth d with two actions, left and right (which traverses down the tree appropriately). Assume the agent receives a reward of 1 if and only if it takes an action leading to a particular fixed leaf. Compute the return of a policy that chooses between left and right uniformly at random.

[Expected Reward] We compute the distribution of states at every timestep, considering that from any given state, the agent has a $1/2$ chance of moving left or right. At each depth k of the tree, the agent can be in 2^k possible states (nodes) since the tree branches out by a factor of 2 at each level. The probability of being in any particular state at depth k is $(1/2)^k$ because to arrive at that state, the agent must make a particular sequence of k choices, each with a probability of $(1/2)$. To compute the expected return, we aggregate the probabilities of reaching the rewarding state over all possible paths. However, since there's only one rewarding path and the reward is located at one leaf, the expected return is $(1/2)^d$. This is because regardless of the path distribution across the tree, only the path leading to the reward contributes to the expected return.

Remark 2.2. Observe that an MDP, together with a fixed policy, induces a particular Markov chain. \lrcorner

Remark 2.3 (On Reward Choice). This example underscores how crucial the choice of reward and state space structure is in reinforcement learning. In a binary tree with a single reward at one leaf

and a uniform random policy, the expected return diminishes exponentially with the depth of the tree. This setup illustrates a near-worst-case scenario for an RL agent, where the environment's structure and reward distribution make learning efficient policies highly challenging without additional strategies or information. It highlights the balance between exploration and exploitation and the importance of reward shaping or structure to guide learning in complex environments. ┘