

**Reinforcement learning
and
Markov Decision Processes (MDPs)**

15-859(B)

Avrim Blum

RL and MDPs

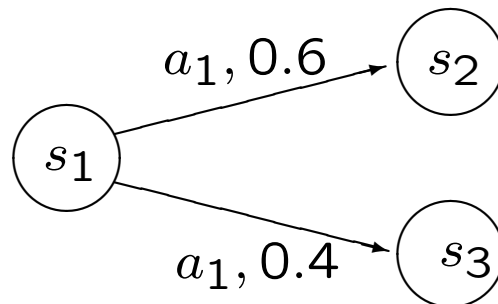
General scenario: We are an **agent** in some **state**. Have observations, perform actions, get rewards. (See lights, pull levers, get cookies)

Markov Decision Process: like DFA problem except we'll assume:

- Transitions are probabilistic. (harder than DFA)
- Observation = state. (easier than DFA)

Assumption is that reward and next state are (probabilistic) functions of current observation and action only.

- Goal is to learn a good strategy for collecting reward, rather than necessarily to make a model. (different from DFA)



Typical example

Imagine a grid world with walls.

- Actions left, right, up, down.
- If not currently hugging a wall, then with some probability the action takes you to an incorrect neighbor.
- Entering top-right corner gives reward of 100 and then takes you to a random state.

Nice features of MDPs

- Like DFA, an appealing model of agent trying to figure out actions to take in the world. Incorporates notion of actions being good or bad for reasons that will only become apparent later.
- Probabilities allow to handle situations like: wanted robot to go forward 1 foot but went forward 1.5 feet instead and turned slightly to right. Or someone randomly picked it up and moved it somewhere else.
- Natural learning algorithms that propagate reward backwards through state space. If get reward 100 in state s , then perhaps give value 90 to state s' you were in right before s .
- Probabilities can to some extent model states that look the same by merging them, though this is not always a great model.

Limitations

- States that look the same can be a real problem. E.g., “door is locked” and “door is unlocked”. Don’t want to just keep trying, and explicitly modeling belief-state blows up problem size.
- Markov assumption not quite right (similar issue).
- “POMDP” model captures probabilistic transitions *and* lack of full observability, but much less to say about them.

What exactly do we mean by a good strategy?

Several notions of what we might want a learned strategy to optimize:

- Expected reward per time step.
- Expected reward in first t steps.
- Expected discounted reward: $r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 + \dots$ ($\gamma < 1$)

We will focus on this last one.

- Why γ^i , and not, e.g., $1/i^2$?

One answer: makes it *time independent*. In other words the best action to take in state s doesn't depend on when you get there.

So, we are looking for an optimal *policy* (a mapping from states to actions).

Q-values

Goal is to maximize discounted reward:

$$r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 + \dots,$$

Define: $Q(s, a)$ = expected discounted reward if perform a from s and then follow optimal policy from then on.

Define: $V(s) = \max_a Q(s, a)$.

Equivalent defn: $V(s) = \max_a [R(s, a) + \gamma \sum_{s'} \text{Pr}(s') V(s')]$.

($R(s, a)$ = expected reward for doing action a in state s .)

Why is this OK as a definition?

A: Only one solution.

Can see this either by proving by contradiction, or noticing that if you're off by some amount Δ on the right-hand-side, then you'll be off by only $\gamma\Delta$ on the left-hand-side.

How to solve for Q-values?

Suppose we are **given** the transition and reward functions. How to solve for Q-values? Two natural ways:

1. **Dynamic Programming.** Start with guesses $V_0(s)$ for all states s . Update using:

$$V_i(s) = \max_a \left[\mathbf{E}[R(s, a)] + \gamma \sum_{s'} \text{Pr}(s') V_{i-1}(s') \right]$$

Get ϵ -close in $O(\frac{1}{\gamma} \log \frac{1}{\epsilon})$ steps. In fact, if initialize all $V_0(s) = 0$, then $V_t(s)$ represents max discounted reward if game ends in t steps.

2. **Linear Programming.** Replace “max” with “ \geq ” and minimize $\sum_s V(s)$ subject to these constraints.

Q-learning

Start off with initial guesses $\hat{Q}(s, a)$ for all Q -values. Then update these as you travel. Update rules:

- Deterministic world: in state s , do a , go to s' , get reward r :

$$\hat{Q}(s, a) \leftarrow r + \gamma \max_{a'} \hat{Q}(s', a')$$

- Probabilistic world: on the t th update of $\hat{Q}(s, a)$:

$$\hat{Q}(s, a) \leftarrow (1 - \alpha_t) \hat{Q}(s, a) + \alpha_t [r + \gamma \max_{a'} \hat{Q}(s', a')]$$

Idea: dampen the randomness.

$\alpha_t = 1/t$ or similar. With $\alpha_t = 1/t$, you get a fair average of all the rewards received for doing a in state s . If you make more slowly decreasing, you favor more recent r 's.

Proof of convergence (deterministic case for simplicity)

Start with some \hat{Q} values. Let: $\Delta_0 = \max_{s,a} |\hat{Q}(s,a) - Q(s,a)|$.

Define a “ Q -interval” to be an interval in which every (s,a) is tried at least once.

Claim: after the i th Q -interval, $\max_{s,a} |\hat{Q}(s,a) - Q(s,a)| \leq \Delta_0 \gamma^i$. In addition, during the i th Q -interval, this maximum difference will be at most $\Delta_0 \gamma^{i-1}$.

Proof: Prove by induction. Base case (the very beginning) is OK. After an update in interval i of $\hat{Q}(s,a)$, (let’s say that action a takes you from s to s') we have:

$$\begin{aligned} |\hat{Q}(s,a) - Q(s,a)| &= |(r + \gamma \max_{a'} \hat{Q}(s',a')) - (r + \gamma \max_{a'} Q(s',a'))| \\ &= \gamma |\max_{a'} \hat{Q}(s',a') - \max_{a'} Q(s',a')| \\ &\leq \gamma \max_{a'} |\hat{Q}(s',a') - Q(s',a')| \\ &\leq \gamma \cdot \gamma^{i-1} \Delta_0. \quad \blacksquare \end{aligned}$$

So, to get convergence, pick actions so that #intervals $\rightarrow \infty$.

Does approximating V give an apx optimal policy?

Say we find \hat{V} and use greedy policy π with respect to \hat{V} . Does $\hat{V}(s) \approx V(s)$ necessarily imply $V^\pi(s) \approx V(s)$ for all states s ?

- $V^\pi(s)$ is value of s if we follow policy π .
- Let $\epsilon = \max_s |\hat{V}(s) - V(s)|$. Assuming this quantity is small.
- Let $\Delta = \max_s [V(s) - V^\pi(s)]$. Want to show this has to be small too.

Does apx V give an apx optimal policy? Yes.

Let s be a state where gap is largest: $V(s) - V^\pi(s) = \Delta$. Say $\text{opt}(s) = a$ but $\pi(s) = b$. $V^\pi(s) = R(s, b) + \gamma \sum_{s'} \text{Pr}_b(s') V^\pi(s')$.

Step 1: Consider following π for 1 step and then doing opt from then on. At best this helps by $\gamma\Delta$.

Step 2: So, this implies that

$$\Delta(1 - \gamma) \leq [R(s, a) + \gamma \sum_{s'} \text{Pr}_a(s') V(s')] - [R(s, b) + \gamma \sum_{s'} \text{Pr}_b(s') V(s')].$$

Step 3: But, since b looked better according to \hat{V} ,

$$0 \geq [R(s, a) + \gamma \sum_{s'} \text{Pr}_a(s') \hat{V}(s')] - [R(s, b) + \gamma \sum_{s'} \text{Pr}_b(s') \hat{V}(s')].$$

Step 4: But since $|V(s') - \hat{V}(s')| \leq \epsilon$, by subtracting we get:

$$\Delta(1 - \gamma) \leq \gamma [\sum_{s'} \text{Pr}_a(s') \epsilon] + \gamma [\sum_{s'} \text{Pr}_b(s') \epsilon] \leq 2\gamma\epsilon.$$

So, $\Delta \leq 2\epsilon\gamma/(1 - \gamma)$.

What if state space is too large to write down explicitly?

In practice, often have large state space. Each state described by set of features.

Much like concept learning except:

- Next example is probabilistic function of action and previous example.
- Most examples don't have labels. Only get feedback infrequently (e.g., when you win the game, reach the goal, etc.).

What if state space is too large to write down explicitly?

Neat idea: Use Q-learning (or TD(λ)) to train standard learning algorithm with “hallucinated” feedback.

Say we're in state s , do action a , get reward r , and go to s' . Use $(1 - \alpha)\hat{Q}(s, a) + \alpha(r + \gamma\hat{V}(s'))$ as the label. Can think of like training up an evaluation function for chess by trying to make it be self-consistent.

Will this really work? “May work in practice but it will never work in theory”.

Depends on how well your predictor can fit the value function. Even if it can fit it, you might still get into a bad feedback loop.

Work by Geoff Gordon on what conditions really ensure things will go well. In practice, it can still sometimes work fine even if these aren't satisfied. E.g., TD-gammon backgammon player.