

## MDPs

## Exercise 1

- (a) **Environment:** Opponent's chess positions (Chess position: a 'instant picture' of the own pieces at particular turn)  
**States:** Player's own pieces' chess position (dim: *num\_of\_diff\_pieces* x *num\_of\_legal\_piece\_pos*, discrete)  
**Actions:** legal moves for all the pieces  $\approx 80$ , discrete  
**Reward:** +1 for winning, -1 for losing, 0 for draw. No intermediate rewards
- (b) **Environment:** 3D continuous space with coordinates of the object to pick up  
**States:** Arm position in 8D space (continuous)  
**Actions:** move the arm in one of the 8 dimensions (continuous) - back and forth, up and down, right and left, rotation  
**Reward:** big positive reward for picking the object, e.g. 100, small negative rewards for each move, e.g. -1 to enforce efficiency.  
The same for the placing procedure...
- (c) **Environment:** 3D continuous space with possibly changing dynamics  
**States:** position of each of the 4 rotors in 3D (continuous)  
**Actions:** thrust power in each of the thrusters (4D, continuous)  
**Reward:** Continuous reward for the time being on the same point on space
- (d) Dialogue agent. Goal: accomplish a task the user asks for: book flight, provide information etc. in a dialogue setting. This is operating directly on utterance level, not on closed set of possible dialogue acts (templates) like most of the actual dialogue agents. The actual task can be accomplished after gathering enough knowledge from the user - turning unstructured text to structured information based on a knowledge graph, and reasoning on it.  
**Environment:** User utterance  
**States:** Discrete information - each user utterance produces new structured input - a edge or vertex on the knowledge graph. It's given from the environment and not agent's task to covert the utterance to structured input. Dimensions depend on the size of the knowledge graph  
**Actions:** Pick a word from vocabulary - 40k (New states are encountered not after each action, but after building an utterance, when the agent hits the EOS token. )  
**Reward:** Small negative for each dialogue turn, big positive reward if user marks task accomplished, big negative reward if human agent has to step in.

**Exercise 2**

1. In the case of MDPs each action is taken given some policy and each action can lead the agent to end in different subsequent state. From these states we can take another action and encounter new rewards, depending on our policy, so we wish to know the future rewards too. The expected return and thus the value of the action-state pair depends also on the policy. In bandit learning we don't maintain states. We end up in the same state after each action, we are not interested in the expectation of sequence of rewards but rather in expectation of the reward after taking this action.

In bandit setting:

$$\begin{aligned} p(s', r|s, a) &= p(r|a), \text{ because } s' = s, |S| = 1 \\ \implies r(s, a) &= \sum_r r \sum_{s'} p(s', r|s, a) = \sum_r p(r|a)r = \mathbb{E}(R_t|A_t = a) \end{aligned}$$

$$\begin{aligned} q(a, s) &= \mathbb{E}(R_t + \gamma R_{t+1} + \dots | S_t = s, A_t = a) \stackrel{\text{bandits}}{=} \\ &= \mathbb{E}(R_t + \gamma R_{t+1} + \dots | A_t = a) + \dots \stackrel{\text{linearity of cond.exp.}}{=} \\ &= \mathbb{E}(R_t | A_t = a) + \mathbb{E}(\gamma R_{t+1} | A_t = a) + \dots \stackrel{\text{bandits}}{=} \\ &= \mathbb{E}(R_t | A_t = a) + \gamma \mathbb{E}(R_t | A_t = a) + \dots \\ \implies &\text{considering future rewards doesn't provide us with no new information} \end{aligned}$$

2. Given:

$$\begin{aligned} \sum_a \pi(a|s) &= 1 \implies \mathbb{E}(G_t | S_t = s, A_t = a) \pi(a|s) = \mathbb{E}(G_t | S_t = s) \\ \implies q_\pi(s, a) \pi(a|s) &= v_\pi(s) \end{aligned}$$

3. From the Bellman equation:

$$\begin{aligned} v_\pi(s) &= \sum_a \pi(a|s) \left( \underbrace{\sum_r \sum_{s'} p(s', r|s, a)r}_{r(s, a)} + \underbrace{\sum_{s'} \sum_r p(s', r|s, a)v_\pi(s')}_{{p(s'|s, a)}} \right) = \\ &= \sum_a \pi(a|s) (\sum_{s'} p(s'|s, a)r(s, a, s') + \sum_{s'} p(s'|s, a)v_\pi(s')) = \\ &= \sum_a \pi(a|s) (\sum_{s'} p(s'|s, a)(r(s, a, s') + v_\pi(s'))) \end{aligned}$$

where  $r(s, a) = \sum_{s'} p(s'|s, a)r(s, a, s')$

### Exercise 3

- (a) Given  $\pi$ : discrete, at each state there can be mostly  $|A|$  different ways of acting. Thus, the number of policies is  $|A|^{|S|}$ . Many of the transitions may have probability 0.
- (b)  $\mathbf{v}_\pi = \mathbf{r} + \gamma \mathbf{P}_\pi \mathbf{v}_\pi$   
 $\mathbf{v}_\pi - \gamma \mathbf{P}_\pi \mathbf{v}_\pi = \mathbf{r}$   
 $\mathbf{v}_\pi (\mathbf{I} - \gamma \mathbf{P}_\pi) = \mathbf{r}$   
 $\mathbf{v}_\pi = (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r}$

Value function for policy\_left (always going left):

```
[0.          0.          0.53691275  0.          0.          1.47651007
 0.          0.          5.          ]
```

Value function for policy\_right (always going right):

```
[0.41401777  0.77456266  1.31147541  0.36398621  0.8185719   2.29508197
 0.13235862  0.          5.          ]
```

As expected, the output shows that always getting right is a better option, leading to the direction of the goal state. Going left can reach the goal only from the rightmost states which are the only ones with value different than 0, otherwise remains away from getting to the goal and thus getting the reward. Also, in both policies there is no escape and thus no probability to get the reward from state 7 and state 8 is the most desired state, being the goal itself.

- (c) Optimal value function:  

```
[0.49756712  0.83213812  1.31147541  0.53617147  0.97690441  2.29508197
 0.3063837   0.          5.          ]
```

number optimal policies: 2

optimal policies:

```
[1 2 2 3 3 2 0 0 0]
[2 2 2 3 3 2 0 0 0]
```

For the case of 9 states as above: Except for the first state that has as optimal actions down or right, all other actions have single optimal action.

- (d) For 16 states: 4194304 policies, much more than for the original setting with 9 states: 16384.  
 The assumption that we can evaluate all policies is not feasible for large numbers of states.