

0 Instructions (Total 40 pts)

Homework is due Tuesday, April 30, 2024 at 23:59pm Central Time. Please refer to <https://courses.grainger.illinois.edu/cs446/sp2024/homework/hw/index.html> for course policy on homeworks and submission instructions.

Version history.

- 1.0. Initial version.

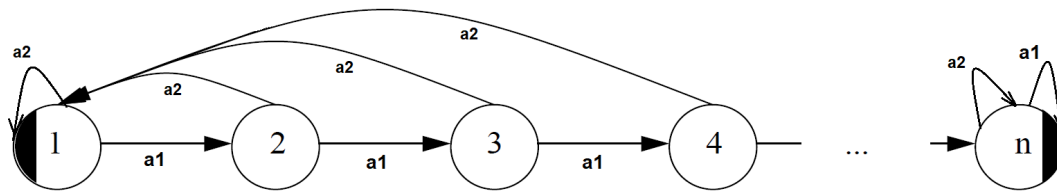
1 Bellman Equation (7 pts)

1. (1 pts) What is the definition of the Q value function, $Q^\pi(s, a)$, for an infinite horizon with discount γ , reward function $R(s, a)$, stochastic transition dynamics $p(s' | s, a)$, with policy $\pi(a | s)$?
2. (2 pts) Derive the Bellman equation from this definition.
3. (1 pts) Given Q and transition (s, a, r, s') , what is the Q-learning update step (for discount γ and learning rate α)?
4. (1 pts) If the maximum reward is $R_{max} = \max_{s,a} R(s, a)$, with discount factor γ , what is the maximum value of the optimal Q, $\max_{s,a} Q^*(s, a)$?
5. (2 pts) Consider a system with two states S_1 and S_2 and two actions a_1 and a_2 . You perform actions and observe the rewards and transitions listed below. Each step lists the current state, reward, action and resulting transition as: $S_i; R = r; a_k : S_i \rightarrow S_j$. Perform Q-learning using a learning rate of $\alpha = 0.5$ and a discount factor of $\gamma = 0.5$ for each step by applying the formula from part (b). The Q-table entries are initialized to zero. Fill in the four tables below corresponding to the following four transitions. What is the final policy after having observed the four transitions?
 - (a) $S_1; R = -10; a_1 : S_1 \rightarrow S_1$
 - (b) $S_1; R = -10; a_2 : S_1 \rightarrow S_2$
 - (c) $S_2; R = 18.5; a_1 : S_2 \rightarrow S_1$
 - (d) $S_1; R = -10; a_2 : S_1 \rightarrow S_2$

2 Combination Lock (7 pts)

Consider an MDP with n states s_1, s_2, \dots, s_n , where there are 2 actions a_1, a_2 at each state. At each state s_i , $i < n$ action a_1 takes the agent to the s_{i+1} and for state s_n , action a_1 is a self-loop which keeps the agent in s_n . At any state, action a_2 takes the agent back to s_1 , except for the last state, which again has a self-loop. See the figure for an overall picture of the setting. $R(s_i, a_j)$ (reward of action a_j at state s_i) is 0, except for (s_n, a_1) , which has a value of 1. The agent takes one action at each step.

With *uniform random policy*, the agent at each state chooses one of the available actions uniformly at random. Now considering this *combination lock* MDP, answer the questions below. You need to show your work to receive full credit for each of the sections.



- (2 pts) Compute the expected number of steps for the uniform random policy to go from state s_1 to state s_n .
- (3 pts) Compute the formula for $Q(s_i, a_j)$, $\forall i, j$ for the uniform random policy considering a discounted reward setting with a discount factor of γ .
- (2 pts) Now consider a new *greedy policy* using the $Q(s, a)$ function you computed in the previous part. Specifically, $\pi(s_i) = \operatorname{argmax}_a Q(s_i, a)$ (i.e., the agent, at each state, chooses the action with the highest value of the Q function computed in part (b)). Compute the new expected number of steps to get from state s_1 to s_k . Hint: how does $Q(s_i, a_1)$ compare to $Q(s_i, a_2)$?

3 Q-value initialization (7 pts)

In this question we will explore the effect of initialization on Q-learning. Let $Q_1(s, a) = 0$, $Q_2(s, a) = R_{max}$, where $R_{max} = \max_{s,a} R(s, a)$. Consider the greedy policy $\pi_i(s) = \operatorname{argmax}_a Q_i(s, a)$, where ties are broken at random (by selecting a random action among equal values).

Consider the MDP in the previous problem:

1. (2 pts) What is the expected number of steps for each policy to reach s_n starting from s_1 ?
2. (1 pts) Now assume that each time we explore a state transition (s, a, s', r) we update Q_i with a Bellman update, i.e.,

$$Q_i(s, a) = (1 - \alpha)Q_i(s, a) + \alpha(r + \gamma \max_{a'} Q_i(s', a'))$$

for $0 < \gamma < 1, 0 < \alpha < 1$.

Now what is the expected number of steps for π_1 to reach s_n starting from s_1 ?

3. (1 pts) Assuming we reset the agent (π_1) to s_1 once it reaches s_n , what is the expected number of steps before it reaches s_n again?
4. (1 pts) How can using a replay buffer help reduce the number of episodes before π_1 is optimal?
5. (2 pts) What are the minimum and the maximum number of steps for π_2 to reach s_n starting from s_1 (again assuming we update Q_2 each step)?

4 Policy Gradient (8 pts)

In policy gradient methods our goal is to evaluate the policy π directly rather than the value function. Let $J(\pi)$ be the expected value of the policy for some initial state distribution d_0 . In practice, π is a parameterized function, say π_θ .

- (3 pts) For trajectory $\tau = s_0, a_0, r_0, \dots, s_T, a_T, r_T$ the discounted sum of rewards is $R(\tau) = \sum_{i=0}^T \gamma^i r_i$ and $P^\pi(\tau)$ is the probability of τ under the policy, i.e., $P^\pi(\tau) = d_0(s_0) \prod_{i=1}^{T-1} \pi(a_i | s_i) P(s_{i+1} | s_i, a_i) \pi(a_T | s_T)$. Show that

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[R(\tau) \sum_{i=0}^T \nabla_\theta \log \pi_\theta(a_i | s_i) \right]$$

- (3 pts) An algorithm using the above gradient requires rolling out trajectories for each update. Let $d_i^\pi(s)$ be the probability the policy visits state s after i steps from the starting state distribution, and let $d^\pi(s)$ be the (discounted) probability that the policy visits s on some (infinite) trajectory, i.e., $d^\pi(s) = \sum_{i=0}^{\infty} \gamma^i d_i^\pi(s)$. Note this is called the stationary distribution, or discounted state occupancy. Show that

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{s \sim d^\pi, a \sim \pi_\theta(s)} [Q^\pi(s, a) \nabla_\theta \log \pi_\theta(a | s)]$$

Hint: Notice that $J(\pi) = \mathbb{E}_{s \sim d_0} [V^\pi(s)]$, where V is the value function, and so $\nabla_\theta J(\pi_\theta) = \nabla_\theta [\mathbb{E}_{s \sim d_0} [V^{\pi_\theta}(s)]]$. Now how do you write $V^{\pi_\theta}(s)$ in terms of Q^{π_θ} ?

- (2 pts) Finally, show that

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{s \sim d^\pi, a \sim \pi_\theta(s)} [(Q^\pi(s, a) - f(s)) \nabla_\theta \log \pi_\theta(a | s)]$$

for any function f . Many types of policy gradient depend on appropriate choice of f to reduce the variance of the estimated gradient.

5 Coding: Tabular Q-Learning (11 pts)

In this problem you will be implementing tabular Q learning for the Taxi-v3 environment.

1. (2 pts) Your first task is to investigate the OpenAI gym API and the Taxi-v3 environment. You can create the environment with

```
env = gym.make("Taxi-v3", render_mode="ansi")
```

- (a) What do states correspond to? How many possible states are there? What is the action space?
 - (b) What do `env.step()` and `env.reset()` return? How are states represented?
 - (c) What are the valid `render_modes` in `gym.make`? For Taxi-v3 what does each render mode do/return?
 - (d) Go digging in the open source documentation - how does `render()` convert the state representation returned by `step` and `reset` (hint: `self.s`) to the human readable features you described in part (a). Your answer should be a function call, `"self.function_name(self.s)"`.
2. (2 pts) Implement the `QAgent` class, more instructions are in the code. Report the `train_rewards.png` plot created by `q_learning()` for an experiment with $\alpha = 0.1, \gamma = 0.9, \epsilon = 0.1$ and 10,000 epochs.
 3. (2 pts) What happens if you initialize the Q table values to -10 (the minimum reward in the environment)? Why does it perform worse? Report the `train_moving_avg.png` plot (with -10 initial Q values) created by `q_learning()` for an experiment with $\alpha = 0.1, \gamma = 0.9, \epsilon = 0.1$ and 10,000 epochs.
 4. (1 pts) What percentage of the time does a uniform random policy ($\epsilon = 1.0$) get positive reward (success) in the environment across 10,000 epochs (printed for you as "Episode success rate")? Report the `train_rewards.png` plot for this uniform random policy.
 5. (1 pts) With initialization of -10 , report the `train_rewards.png` plot created by `q_learning()` for an experiment with $\alpha = 0.1, \gamma = 0.9$, 10000 epochs and ϵ of your choice which yields a model that successfully passes the evaluation (meaning the taxi delivers the person to the goal in the call to `"eval_agent()"`).
 6. (2 pts) Some actions lead to "no-op", for example moving right when there is a wall to the right. In the documentation for the Taxi-v3 environment they describe a method to only sample actions which lead to a change in state. Implement this; now, what percentage of the time does a uniform random policy ($\epsilon = 1.0$) get positive reward (success) in the environment across 10,000 epochs? Report the `train_rewards.png` plot.
 7. (1 pts) Again, with initialization of -10 AND never taking an action that leads to "no-op", report the `train_rewards.png` plot created by `q_learning()` for an experiment with $\alpha = 0.1, \gamma = 0.9$, 10000 epochs and ϵ the same as in part 5 which yields a model that successfully passes the evaluation.