

1. Learning by Example

1.1. $V^\pi(s) = 0$ for any state

$$\text{Episode 1: } V^\pi(A) = (1-0.25) * V^\pi(A) + 0.25[R(A, \text{right}, B) + V^\pi(B)]$$

$$V^\pi(A) = (0.75) * 0 + 0.25[16 + 0] = 4$$

$$\text{Episode 2: } V^\pi(B) = (1-0.25) * V^\pi(B) + 0.25[R(B, \text{right}, C) + V^\pi(C)]$$

$$V^\pi(B) = (0.75) * 0 + 0.25[4 + 0] = 1$$

$$\text{Episode 3: } V^\pi(B) = (1-0.25) * V^\pi(B) + 0.25[R(B, \text{down}, E) + V^\pi(E)]$$

$$V^\pi(B) = (0.75) * 1 + 0.25[-12 + 0] = -2.25$$

$$\text{Episode 4: } V^\pi(C) = (1-0.25) * V^\pi(C) + 0.25[R(C, \text{down}, F) + V^\pi(F)]$$

$$V^\pi(C) = (0.75) * 0 + 0.25[-16 + 0] = -4$$

$$\text{Episode 5: } V^\pi(F) = (1-0.25) * V^\pi(F) + 0.25[R(F, \text{stay}, F) + V^\pi(F)]$$

$$V^\pi(F) = (0.75) * 0 + 0.25[8 + 0] = 2$$

$$\text{Episode 6: } V^\pi(C) = (1-0.25) * V^\pi(C) + 0.25[R(C, \text{down}, F) + V^\pi(F)]$$

$$V^\pi(C) = (0.75) * -4 + 0.25[-9 + 2] = -4.75$$

State	$U^\pi(\text{state})$
A	4
B	-2.25
C	-4.75
D	0
E	0
F	2

2. Reinforcements

$$2.1. \quad Q(A, \text{right}) = (1-0.5) * Q(A, \text{right}) + 0.5(r_1 + \max Q(s_B, a'))$$

$$= 0.5 * 0 + 0.5(4 + 0) = 2$$

$$Q(C, \text{left}) = (1-0.5) * Q(C, \text{left}) + 0.5(r_1 + \max Q(s_B, a'))$$

$$= 0.5 * 0 + 0.5(4 + 0) = 2$$

$$Q(B, \text{right}) = (1-0.5) * Q(B, \text{right}) + 0.5(r_1 + \max Q(s_C, a'))$$

$$= 0.5 * 0 + 0.5(-4 + 2) = -1$$

$$Q(A, \text{right}) = (1-0.5) * Q(A, \text{right}) + 0.5(r_1 + \max Q(s_B, a'))$$

$$= 0.5 * 2 + 0.5(8 + 0) = 5$$

2.1.1. $Q(A, \text{right}) = 5$
 $Q(B, \text{right}) = -1$

2.1.2. $\pi_Q(A) = \text{right}$
 $\pi_Q(B) = \text{left}$

2.2. $T(A, \text{right}, B) = 1.0$ (was at state A two times and went right to state B twice)
 $R(A, \text{right}, B) = (4 + 8) / 2 = 6$

$T(B, \text{right}, A) = 0$ (was at state B once and went right and the next state was state C, not state A)

$R(B, \text{right}, A) = \text{N/A}$ (since there was no observed transition from B, going right, to A, there was no reward)

$T(B, \text{left}, A) = 0$ (was at state B once, but went right to state C)

$R(B, \text{left}, A) = \text{N/A}$ (since there was no observed transition from B, going left, to A, there was no reward)

3. Reinforcement Learning Background

3.1. True or false

i) False, Temporal difference learning updates its value estimates incrementally, made from interactions with the environment

ii) True, while optimal, exploration involves uncertainty which, in some cases, can lead to regret due to suboptimal decisions (from exploring new actions).

iii) True, since there is no randomness in a deterministic MDP, with a learning rate of $\alpha = 1$, it will eventually converge to the optimal policy

iv) False, having a large discount (close to 1) means that the agent values the future rewards almost as much as the immediate rewards. Greedy behavior, on the other hand, places an emphasis on immediate rewards and not the future rewards.

v) True, a large, negative living reward incentivizes the agent to avoid prolonged interaction with the environment because it keeps losing rewards over time. Thus, the agent is encouraged to prioritize immediate rewards, which is a form of greedy behavior

vi) False, a negative living reward isn't always expressed by using a discount < 1 . The discount value impacts how future rewards are valued. It doesn't penalize the agent for staying in the environment, it makes the agent less motivated to

consider future rewards. A negative living reward is conceptually different from the discount value.

vii) True, similar to above, a negative living reward is different from a discount value < 1 , thus it cannot always be expressed as the same.

3.2.

a)

Direct Evaluation to estimate $U(s)$

Temporal Difference learning to estimate $U(s)$

b)

A fixed policy taking actions uniformly at random

An ϵ -greedy policy