1) Which of the following expressions is the TD error $\delta_t$ at time step $t$ in the $TD(0)$ algorithm? All the symbols have their usual meaning.   **1 point**

○ $V(S_{t+1}) - V(S_t)$

○ $R_{t+1} + \gamma \cdot V(S_t) - V(S_t)$

◉ $R_{t+1} + \gamma \cdot V(S_{t+1}) - V(S_t)$   *(selected)*

○ $R_{t+1} + \gamma \cdot V(S_{t+1})$

$$\delta_t \stackrel{\alpha}{=} \left[ R_{t+1} + \gamma\, V(S_{t+1}, a) - V(S_t) \right]$$

TD error

2) Is the following statement true or false?   **1 point**

TD methods require the complete knowledge of a model of the environment.

○ True

◉ False

no, it works on sampling
MC → sampling
DP → bootstrapping

An MDP has two non-terminal states, $A$ and $B$ and a terminal state $C$. Consider the following trajectory under some policy $\pi$:

$A \xrightarrow{2} B \xrightarrow{1} A \xrightarrow{-1} C$

The rewards are mentioned above the transitions. $TD(0)$ is used to evaluate the value function for this policy. The values for $A$ and $B$ are initialized to zero. Use $\alpha = 0.1$ and $\gamma = 1$.

$C \xleftarrow{-1} A \xrightarrow{2} B$

3) What is the value of $V(A)$ at the end of this trajectory?   **1 point**

0.08

$V(A) = V(A) + 0.1(2 + 0 + 0)$
$= 0.2$

4) What is the value of $V(B)$ at the end of this trajectory?   **1 point**

0.12

$V(B) = V(B) + 0.1(1 + 0.2 + 0)$
$= 0.1(1 \cdot 2) = 0.12$

5) In batch-$TD(0)$, we treat all available experience as a batch, compute the sum of all the increments (TD-errors) for the batch, and update the values just once at the end using this sum. This forms one iteration of batch-TD. We do this repeatedly until the value function converges.

If this trajectory is repeatedly used to update the values using batch-$TD(0)$ what does $V(A)$ converge to?   **1 point**

2     $-1 + 1 + 2 = 2$

$V(A) = 0.2 + 0.1(-1 + 0 - 0.2)$
$= 0.2 + 0.1(-1.2)$
$= 0.2 + (-0.12)$
$= 0.08$

$V(A) = V(A) + \alpha \left[ R + \gamma V(B) - V(A) \right]$

$V(A) = \frac{1}{3}\left[ 2 + (1 + V(A)) \right]$
$+ \frac{1}{3}\left[ 1 + (1 + V(B)) \right]$
$+ \frac{1}{3}(-1)$
$\downarrow V(C)$

6) In batch-$TD(0)$, we treat all available experience as a batch, compute the sum of all the increments (TD-errors) for the batch, and update the values just once at the end using this sum. This forms one iteration of batch-TD. We do this repeatedly until the value function converges.

If this trajectory is repeatedly used to update the values using batch-$TD(0)$ what does $V(B)$ converge to?

3

*2+1 ⇒ reward func'n estimate*

7) Which of the following conditions are necessary for ensuring the convergence of SARSA? **1 point**

☐ Each state-action pair is visited at least once.

☑ All state-action pairs are visited an infinite number of times.

☑ The policy converges in the limit to a greedy policy.

*① ∞ no. times visit $(s,a)$ pairs*

*② limit to the greedy policy*

☐ The value of $\epsilon$ is fixed to some small value throughout the algorithm.

8) Choose the correct qualifiers for Q-learning from the options given below. **1 point**

☐ On-policy

*off-policy, TD-control*

*SARSA ⇒ on-policy, TD-control*

☑ Off-policy

☑ TD-control

☐ TD-prediction

9) Consider this statement "The action value function corresponding to the optimal policy is learnt while the actions are sampled from an arbitrary policy." **1 point**
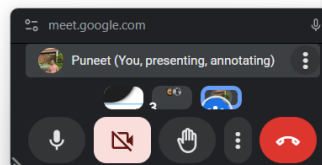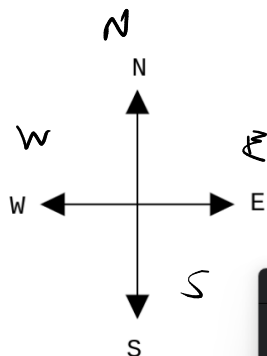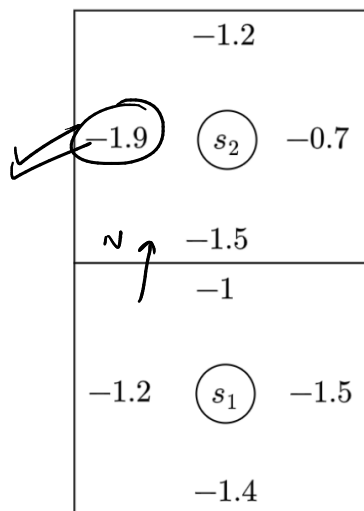
This is true for which of the following algorithms?

○ $TD(0)$

○ SARSA

◉ Q-Learning

Consider a grid world with deterministic transitions and a unit reward of $-1$ for all time-steps. SARSA is run on this setup. The current estimate for the action values for certain states are displayed in the figure.

```
        -1.2
  -1.9  (s₂)  -0.7
  N    -1.5
       -1
  -1.2  (s₁)  -1.5
        -1.4
```

```
       N
       ↑
   W   |
 W ←---+---→ E
       |
       ↓   S
       S
```

meet.google.com
Puneet (You, presenting, annotating)

The agent is currently in state $s_1$. The policy is $\epsilon$-greedy ($\epsilon = 0.1$) with respect to the current estimate of the action values. The action to be executed at this time step in the episode is north. The action that the agent has committed to take from the next state (in this episode) turns out to be the worst possible non-greedy action from that state.

10) Perform one update of the action value for the pair $(s_1, \text{north})$. Use $\alpha = 0.1, \gamma = 1$. Enter the exact numerical answer. **1 point**

-1.19

$$q(s,a) = q(s,a) + \alpha[a_{t+1} + \gamma\, q(s_{t+1}, a_{t+1}) - q(s,a)]$$

$$q(s_1, \text{north}) = -1 \qquad \Rightarrow \quad q(s_2, a') = -1.9$$

$$= -1 + 0.1(-1 + (-1.9) - (-1))$$

$$= -1 - 0.19 = -1.19$$

11) What will be the next action taken by the agent in this episode?

**1 point**

- ● The action west with a probability of 1 ✓ ✓ *(it is clear from the question itself.*
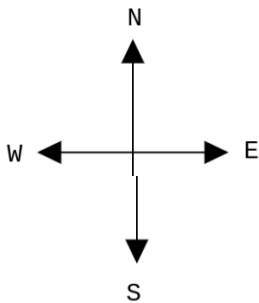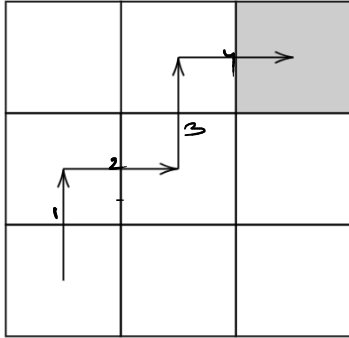- ○ The action west with a probability of 0.1
- ○ The action east with a probability of 0.925 and any of the other three actions with a probability of 0.025
- ○ The action east with a probability of 0.9 and any of the other three actions with a probability of 0.1

12) Compute the importance sampling ratio for the trajectory shown in the figure if the behaviour policy $\mu$ is the equiprobable random policy and the target policy is $\pi$ is as shown in the figure (only for the relevant states):

## Trajectory



```
        N
        ↑
W ◄─────┼─────► E
        ↓
        S
```

$\pi$  Target Policy



Enter your answer correct to two decimal places

186.6

**1 point**

$$W = \frac{\overset{\text{Target}}{\pi(a_1|s_1) \times \pi(a_2|s_2) \cdots}}{\underset{\text{behavioural}}{\mu(a_1|s_1) \times \mu(a_2|s_2) \cdots}} = \frac{0.9 \times 0.9 \times 0.9 \times 1}{(0.25)^4} = 0.729$$

$$= \frac{0.729}{(0.0625)^2} = 186.6$$

## Graded Assignment

1) Match the methods with their corresponding characteristics.

DP: bootstraps, full backups

○ MC: does not bootstrap, full backups

TD: bootstraps, full backups
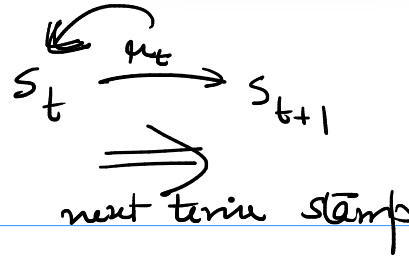
DP: bootstraps, full backups

✓ ○ MC: does not bootstrap, sample backups

TD: bootstraps, sample backups

**2)** What does the quantity $R_{t+1} + \gamma \cdot V(S_{t+1})$ represent in the $TD(0)$ algorithm? All the symbols have their usual meaning. **1 point**

◉ It is a prediction of the expected return of state $S_t$ at time $t$

☑ It is a prediction of the expected return of state $S_t$ at time $t+1$

○ It is a prediction of the expected return of state $S_{t+1}$ at time $t$

○ It is a prediction of the expected return of state $S_{t+1}$ at time $t+1$

$$S_t \xrightarrow{\mu_t} S_{t+1}$$

next term stamp

**3)** Consider a grid world that permits all four actions to be taken from each state. Which of the following are $\epsilon$-soft policies, with $\epsilon = 0.1$? $\mathcal{A} = \{\text{north}, \text{south}, \text{east}, \text{west}\}$. **1 point**

☐ $\pi(a \mid s) = \begin{cases} 1, & a = \text{north} \\ 0, & \text{otherwise} \end{cases}$

☑ $\pi(a \mid s) = 1/4, \quad \forall a \in \mathcal{A}$

☐ $\pi(a \mid s) = \begin{cases} 0.99, & a = \text{north} \\ 0.01/3, & \text{otherwise} \end{cases}$

☑ $\pi(a \mid s) = \begin{cases} 0.925, & a = \text{north} \\ 0.025, & \text{otherwise} \end{cases}$

$\mathcal{E} = 0.1$

$A = \{N, S, E, W\}$

we need probab. of each ac$^n$

$\geq \dfrac{\mathcal{E}}{4} = \dfrac{0.1}{4} = 0.025$

**4)** Choose the correct qualifiers for SARSA from the options given below. **1 point**

☑ On-policy

☐ Off-policy

☑ TD-control

☐ TD-prediction

**5)** Consider the two expressions given below: **1 point**

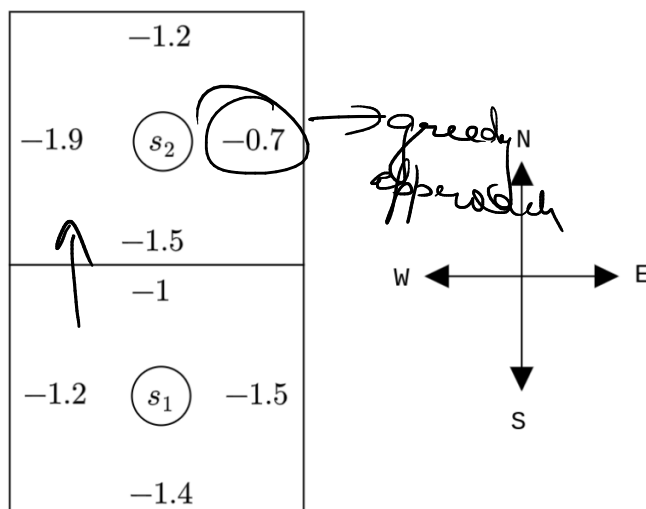$E_1$:
$R_{t+1} + \gamma \cdot Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)$ → SARSA
$E_2$:
$R_{t+1} + \gamma \cdot \max_a Q(S_{t+1}, a) - Q(S_t, A_t)$ → max → Q learning

☑ $E_1$ is the TD error that is used to update the action value function in SARSA

☐ $E_1$ is the TD error that is used to updated the action value function in Q-learning

☐ $E_2$ is the TD error that is used to update the action value function in SARSA

☑ $E_2$ is the TD error that is used to updated the action value function in Q-learning

Consider a grid world with deterministic transitions and a unit reward of $-1$ for all time-steps. Q-learning is run on this setup. The current estimate for the action values for certain states are displayed in the figure.

greedy approach

N, W, E, S (compass)

$-1.2$
$-1.9$ $s_2$ $-0.7$
$-1.5$
$-1$
$-1.2$ $s_1$ $-1.5$
$-1.4$

The agent is currently in state $s_1$. The behaviour policy is $\epsilon$-greedy $(\epsilon = 0.1)$ with respect to the current estimate of the action values. The action $a$

6) Perform one update of the action value for the pair $(s_1, a)$. Use $\alpha = 0.1, \gamma = 1$. Enter the exact numerical answer.

greedy wrt. rewards

$$q(s_1, a) = q(s_1, a) + \alpha \left[ u_{+1} + \gamma \max_a q(s_{t+1}, a) - q(s, a) \right]$$

$$= -1 + 0.1 \left[ -1 + (-0.7) - (-1) \right]$$

$$= -1 + (-0.07) = -1.07$$

7) What will be the next action taken by the agent in this episode?                    **1 point**

○ The action east with a probability of $1$

○ The action east with a probability of $0.1$

✓ The action east with a probability of $0.925$ and one of the other three actions with a probability of $0.025$ each

4 possible states          $\varepsilon = 0.1$

$$\frac{0.1}{4} = 0.025$$

○ The action east with a probability of $0.9$ and one of the other three actions with a probability of $0.1$ each

Consider an undiscounted, episodic task that has two non-terminal states, $A, B$. The rewards are binary, either $1$ or $0$. The following are some episodes experienced by an agent following a fixed policy. The terminal state is not explicitly mentioned for any of the episodes.

$A, 1, B, 0$  $= 1$
$A, 0, B, 1$  $= 1$
$A, 1, B, 0$  $= 1$
$A, 1, B, 0$  $= 1$
$B, 1$  $= 1$
$B, 0$  $= 0$
$B, 1$  $= 1$
$B, 1$  $= 0$

8) What is the estimate of $V(A)$ returned by first-visit MC?

1

$$MC \Rightarrow V(A) = \frac{\text{total reward}}{\text{total traject}} = \frac{4}{4} = 1$$

9) What is the estimate of $V(B)$ returned by first-visit MC?

0.5

$$V(B) = \frac{4}{8} = 0.5$$

Yes, the answer is correct.
Score: 1

Accepted Answers:

(Type: Numeric) 0.5

$\to$ undiscounted $= 1$

10) What is the estimate of $V(A)$ returned by batch-TD(0)?    $$V(A) = E\left[ x + \gamma V(B) \right]$$

1.25

Yes, the answer is correct.
Score: 1

Accepted Answers:

(Type: Numeric) 1.25

$1 + 0.5 = 1.5$  } 2
$0 + 0.5 = 0.5$
$1 + 0.5 = 1.5$  } 3
$1 + 0.5 = 1.5$

11) What is the estimate of $V(B)$ returned by batch-TD(0)?    $$V(B) = 0.5$$    $\frac{5}{4} = 1.25$

0.625  ✗  $0.5$

because B after is terminal

12) Which of the following is the MDP corresponding to the certainty equivalence estimate for this data? The transition probability and the expected immediate reward for a transition are $p, r$ and it is this quantity that is written over the edges.

*(handwritten: TD —)*

*(handwritten top right: MC → least sq. estimate; TD → certainity equivalent)*

**(1)**

- 0.5, 1 → B
- 0.5, 1
- A
- 0.5, 1
- terminal

*(circled (1))*

**(2)**

- 1, 1 → B
- 1, 1
- A
- terminal

**(3)**

- 1, 0.75 → B
- $\frac{3}{4}$
- 1, 0.5
- A
- terminal

*(handwritten: B $\frac{4}{8}$ (0.5) → ½)*

*(check mark by (3))*

Enter 1, 2 or 3 as your answer.

13) Consider the off-policy MC method for evaluating a target policy $\pi$. If the behaviour policy is $\mu$, which of the following statements is an expression of the assumption of coverage?     **1 point**

*(handwritten: $\pi(a|s) > 0 \implies \mu(a|s) > 0$)*

- $\pi(a \mid s) > 0$ whenever $\mu(a \mid s) > 0$
- $\mu(a \mid s) > 0$ whenever $\pi(a \mid s) > 0$
- $\mu(a \mid s)$ should always be greater than $0$ and this is independent of $\pi$
- $\mu(a \mid s) = \pi(a \mid s)$