# AI 3000 / CS 5500 : Reinforcement Learning
## Assignment № 1

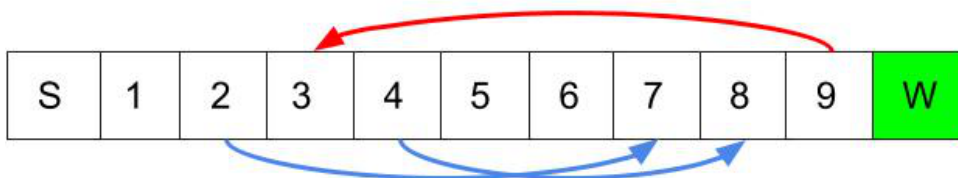**Due Date : 01/09/2022**

---

Course Instructor : Easwar Subramanian                                            21/08/2022

## Problem 1 : Markov Reward Process

Consider the following snake and ladders game as depicted in the figure below.



- Initial state is $S$ and a fair four sided die is used to decide the next state at each time

- Player must land exactly on state $W$ to win

- Die throws that take you further than state $W$ leave the state unchanged

(a) Identify the states, transition matrix of this Markov process.                    (1 point)

(b) Construct a suitable reward function, discount factor and use the Bellman equation for the Markov reward process to compute how long does it take "on average" (the expected number of die throws) to reach the state $W$ from any other state.                    (4 points)

## Problem 2 : Markov Decision Process

A production facility has $N$ machines. If a machine starts up correctly in the morning, it renders a daily revenue of $1\$$. A machine that does not start up correctly, needs to be repaired. A visit by a repair man costs $\frac{N}{2}\$$ per day and he repairs all broken machines on the same day. The repair cost is a lump-sum amount and does not depend on the number of machines that is repaired. A machine that has been repaired always starts up correctly the next day. The number of machines that start up correctly the next day depends on the number of properly working machines at present day and is governed by the probability distribution given in the table below, where $m$ stands for the number of (presently) working machines and $n$ stands for the number of ones that would start up correctly the next day. The goal for the facility manager

is to maximize the profits (revenue - costs) earned.

| $m$ | $n=0$ | $n=1$ | $n=2$ | $n=3$ | $\cdots$ | $n=N-1$ | $n=N$ |
|---|---|---|---|---|---|---|---|
| $m=1$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $0$ | $0$ | $0$ | $0$ | $0$ |
| $m=2$ | $\frac{1}{3}$ | $\frac{1}{3}$ | $\frac{1}{3}$ | $0$ | $0$ | $0$ | $0$ |
| $m=3$ | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{4}$ | $0$ | $0$ | $0$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $m=N-1$ | $\frac{1}{N}$ | $\frac{1}{N}$ | $\frac{1}{N}$ | $\frac{1}{N}$ | $\frac{1}{N}$ | $\frac{1}{N}$ | $0$ |
| $m=N$ | $\frac{1}{N+1}$ | $\frac{1}{N+1}$ | $\frac{1}{N+1}$ | $\frac{1}{N+1}$ | $\frac{1}{N+1}$ | $\frac{1}{N+1}$ | $\frac{1}{N+1}$ |

(a) Formulate the above problem as a Markov decision process by ennumerating the state space, action space, rewards and transition probabilities. (3 Points)

(b) Would you use discounted or undiscounted setting for the above MDP formulation ? Justify the answer. (1 Point)

(c) Suppose the facility manager adopts the policy to never call the repair man. Calculate the value of the policy. For this sub-problem assume that the number of machines in the facility to be five. (3 Points)

(d) Perform one iteration the of policy iteration algorithm on the no-repair policy adopted by the facility manager to get an improved policy for the five machine scenario. (3 Points)

## Problem 3 : On Ordering of Policies

Consider the MDP shown in Figure 1. The MDP has 4 states $\mathcal{S} = \{A, B, C, D\}$ and there are two actions $a_1$ and $a_2$ possible. The actions determine which direction to move from a given state. We consider a stochastic environment such that action suggested by the policy succeeds 90 % of the times and fails 10 % of the times. Upon failure, the agent moves in the direction suggested by the other action. The state $D$ is a terminal state with reward of 100. One can think that terminal states have only one action (an exit option) which gives the terminal reward 100. We consider three policies to this MDP.

- Policy $\pi_1$ is deterministic policy that chooses action $a_1$ at all states $s \in \mathcal{S}$.

- Policy $\pi_2$ is another deterministic policy that chooses action $a_2$ at all states $s \in \mathcal{S}$.

- Policy $\pi_3$ is a stochastic policy described as follows

    - Action $a_1$ is chosen in states $B$ and $D$ with probability 1.0
    - Action $a_2$ is chosen in state $C$ with probability 1.0
    - Action $a_1$ is chosen in state $A$ with probability $0.4$ and action $a_2$ is chosen with probability $0.6$

(a) Evaluate $V^{\pi}(s)$ for each policy described above using the Bellman evaluation equation for all states $s \in \mathcal{S}$. (3 Points)
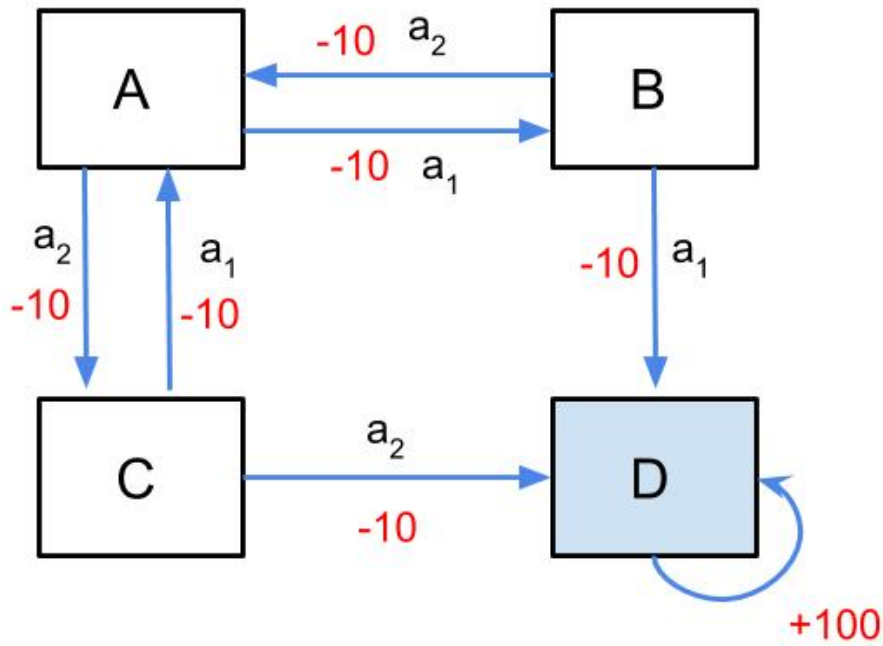
Figure 1: Partial Ordering of Policies

(b) Which is the best policy among the suggested policies ? Why ? (1 Point)

(c) Are all policies comparable ? Provide reason for your answer. (1 Point)

(d) Let $\pi_1$ and $\pi_2$ be two deterministic stationary policies of an MDP $M$. Construct a new policy $\pi$ that is better than policies $\pi_1$ and $\pi_2$. Explain the answer. (3 Points)

[**Note** : $M$ in sub-question (d) is any arbitrary MDP]

## Problem 4 : Effect of Noise and Discounting

Consider the grid world problem shown in Figure 2. The grid has two terminal states with positive payoff (+1 and +10). The bottom row is a cliff where each state is a terminal state with negative payoff (-10). The greyed squares in the grid are walls. The agent starts from the yellow state $S$. As usual, the agent has four actions $\mathcal{A} = $ (Left, Right, Up, Down) to choose from any non-terminal state and the actions that take the agent off the grid leaves the state unchanged. Notice that, if agent follows the dashed path, it needs to be careful not to step into any terminal state at the bottom row that has negative payoff. There are four possible (optimal) paths that an agent can take.

- Prefer the close exit (state with reward +1) but risk the cliff (dashed path to +1)

- Prefer the distant exit (state with reward +10) but risk the cliff (dashed path to +10)

- Prefer the close exit (state with reward +1) by avoiding the cliff (solid path to +1)
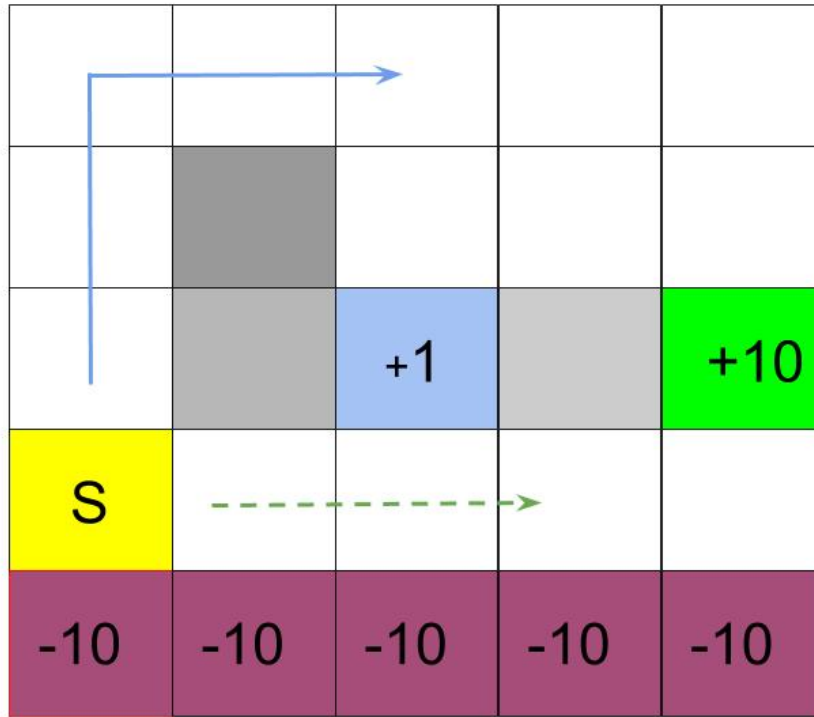
Figure 2: Modified Grid World

- Prefer the distant exit (state with reward +10) by avoiding the cliff (solid path to +10)

There are two free parameters to this problem. One is the discount factor $\gamma$ and the other is the noise factor ($\eta$) in the environment. Noise makes the environment stochastic. For example, a noise of 0.2 would mean the action of the agent is successful only 80 % of the times. The rest 20 % of the time, the agent may end up in an unintended state after having chosen an action.

(a) Identify what values of $\gamma$ and $\eta$ lead to each of the optimal paths listed above with reasoning. (8 Points)

[**Hint** : For the discount factor, consider high and low $\gamma$ values like 0.9 and 0.1 respectively. For noise, consider deterministic and stochastic environment with noise level $\eta$ being 0 or 0.5 respectively]

## Problem 5 : Value Functions

Let $M_1$ and $M_2$ be two identical MDPs with $|\mathcal{S}| < \infty$ and $|\mathcal{A}| < \infty$ except for reward formulation. That is, $M_1 = <\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}_1, \gamma>$ and $M_2 = <\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}_2, \gamma>$. Let $M_3$ be another MDP such that $M_3 = <\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}_1 + \mathcal{R}_2, \gamma>$. Assume the discount factor $\gamma$ to be less than 1.

(a) For an arbitrary but fixed policy $\pi$, suppose we are given action value functions $Q_1^\pi(s, a)$ and $Q_2^\pi(s, a)$, corresponding to MDPs $M_1$ and $M_2$, respectively. Explain whether it is possible to combine these action value functions in a simple manner to calculate $Q_3^\pi(s, a)$ corresponding to MDP $M_3$. (2 Points)

(b) Suppose we are given optimal polices $\pi_1^*$ and $\pi_2^*$ corresponding to MDPs $M_1$ and $M_2$, respectively. Explain whether it is possible to combine these optimal policies in a simple manner to formulate an optimal policy $\pi_3^*$ corresponding to MDP $M_3$. (2 Points)

(c) Suppose $\pi^*$ is an optimal policy for both MDPs $M_1$ and $M_2$. Will $\pi^*$ also be an optimal policy for MDP $M_3$ ? Justify the answer. (2 Points)

(d) Let $\varepsilon$ be a fixed constant. Assume that the reward functions $\mathcal{R}_1$ and $\mathcal{R}_2$ are related as

$$\mathcal{R}_1(s, a, s') - \mathcal{R}_2(s, a, s') = \varepsilon$$

for all $s, s' \in \mathcal{S}$ and $a \in \mathcal{A}$. Let $\pi$ be an arbitrary policy and let $V_1^\pi(s)$ and $V_2^\pi(s)$ be the corresponding value functions of policy $\pi$ for MDPs $M_1$ and $M_2$, respectively. Derive an expression that relates $V_1^\pi(s)$ to $V_2^\pi(s)$ for all $s \in \mathcal{S}$. (3 Points)

# ALL THE BEST