# Assignment #2

Course: *Reinforcement Learning (CS6700)*

Instructor: *Prashanth L.A.*

Due date: *April 20th, 2021*

**Question 1.** Consider a problem of a taxi driver, who serves three cities A, B and C. The taxi driver can find a new ride by choosing one of the following actions.

1. Cruise the streets looking for a passenger.

2. Go to the nearest taxi stand and wait in line.

3. Wait for a call from the dispatcher (this is not possible in town B because of poor reception).

For a given town and a given action, there is a probability that the next trip will go to each of the towns A, B and C and a corresponding reward in monetary units associated with each such trip. This reward represents the income from the trip after all necessary expenses have been deducted. Please refer Table 1 below for the rewards and transition probabilities. In Table 1 below, $p_{ij}^k$ is the probability of getting a ride to town $j$, by choosing an action $k$ while the driver was in town $i$ and $r_{ij}^k$ is the immediate reward of getting a ride to town $j$, by choosing an action $k$ while the driver was in town $i$.

| Town $i$ | Actions $k$ | Probabilities $p_{ij}^k$ $j =$ A  B  C | Rewards $r_{ij}^k$ $j =$ A B C |
|---|---|---|---|
| A | 1 2 3 | $\begin{array}{ccc} A & B & C \\ 1/2 & 1/4 & 1/4 \\ 1/16 & 3/4 & 3/16 \\ 1/4 & 1/8 & 5/8 \end{array}$ | $\begin{array}{ccc} A & B & C \\ 10 & 4 & 8 \\ 8 & 2 & 4 \\ 4 & 6 & 4 \end{array}$ |
| B | 1 2 | $\begin{array}{ccc} A & B & C \\ 1/2 & 0 & 1/2 \\ 1/16 & 7/8 & 1/16 \end{array}$ | $\begin{array}{ccc} A & B & C \\ 14 & 0 & 18 \\ 8 & 16 & 8 \end{array}$ |
| C | 1 2 3 | $\begin{array}{ccc} A & B & C \\ 1/4 & 1/4 & 1/2 \\ 1/8 & 3/4 & 1/8 \\ 3/4 & 1/16 & 3/16 \end{array}$ | $\begin{array}{ccc} A & B & C \\ 10 & 2 & 8 \\ 6 & 4 & 2 \\ 4 & 0 & 8 \end{array}$ |

Table 1: Taxi Problem: Probabilities and Rewards

Suppose $1 - \gamma$ is the probability that the taxi will breakdown before the next trip. The driver's goal is to maximize the total reward untill his taxi breakdown.

Implement the following.             $(1.5 + 1 + 1.5 + 1.5 + 1.5 \text{ marks})$

**1.1**: Find an optimal policy using **policy iteration**(Algorithm 3) starting with a policy that will always cruise independent of the town, and a zero value vector. Let $\gamma = 0.9$.

**1.2**: Run **policy iteration** for discount factors $\gamma$ ranging from 0 to 0.95 with intervals of 0.05 and display the results.

**1.3**: Find an optimal policy using **modified policy iteration**(Algorithm 4) starting with a policy that will always cruise independent of the town, and a zero value vector. Let $\gamma = 0.9$ and $m = 5$.

**1.4**: Find optimal values using **value iteration**(Algorithm 1) starting with a zero vector. Let $\gamma = 0.9$.

**1.5**: Find optimal values using **Gauss-Seidel value iteration**(Algorithm 2) starting with a zero vector. Let $\gamma = 0.9$.

Answer the following questions.             $(1 + 1 + 1 \text{ marks})$

**1.a** How is different values of $\gamma$ affecting the **policy iteration** from **1.2**? Explain your findings.

**1.b** For **modified policy iteration** from **1.3**, do you find any improvement if you choose $m = 10$? Explain your findings.

**1.c** Compare and contrast the behavior of **value iteration** from **1.4** and **Gauss-Seidel value iteration** from **1.5**.

The pseudocode for the Algorithms are given below. (courtesy: Sutton & Barto 1998)

---

**Algorithm 1** Value Iteration

---

1: **Initialize**: $J(s) = 0, \forall s \in \mathcal{S}$;
2: **repeat**
3:      $\delta = 0$;
4:      **for** each $s \in \mathcal{S}$ **do**
5:          $H(s) = \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} P_{ss'}(a)[r(s, a, s') + \gamma J(s')]$
6:          $\delta = \max(\delta, |J(s) - H(s)|)$;
7:      **end for**
8:      **for** each $s \in \mathcal{S}$ **do**
9:          $J(s) = H(s)$;
10:      **end for**
11: **until** $(\delta < 1\mathrm{e}{-8})$

---

---

**Algorithm 2** Gauss-Seidel Value Iteration

---

1: **Initialize**: $J(s) = 0, \forall s \in \mathcal{S}$;
2: **repeat**
3:     $\delta = 0$;
4:     **for** each $s \in \mathcal{S}$ **do**
5:       $j = J(s)$;
6:       $J(s) = \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} P_{ss'}(a)[r(s, a, s') + \gamma J(s')]$
7:       $\delta = \max(\delta, |j - J(s)|)$;
8:     **end for**
9: **until** $(\delta < 1\mathrm{e}{-8})$

---

**Algorithm 3** Policy Iteration

---

1: **Input**: $\pi_0(s), \forall s \in \mathcal{S}$;
2: **Initialize**: $J(s) = 0, \pi(s) = \pi_0(s), \forall s \in \mathcal{S}$;
3: **repeat**
4:     **repeat**
5:       $\delta = 0$;
6:       **for** each $s \in \mathcal{S}$ **do**
7:         $j = J(s)$;
8:         $J(s) = \sum_{s' \in \mathcal{S}} P_{ss'}(\pi(s))[r(s, \pi(s), s') + \gamma J(s')]$
9:         $\delta = \max(\delta, |j - J(s)|)$;
10:      **end for**
11:     **until** $(\delta < 1\mathrm{e}{-8})$
12:     $done = 1$;
13:     **for** each $s \in \mathcal{S}$ **do**
14:       $b = \pi(s)$;
15:       $\pi(s) = \mathrm{argmax}_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} P_{ss'}(a)[r(s, a, s') + \gamma J(s')]$;
16:       **if** $b \neq \pi(s)$ **then**
17:         $done = 0$;
18:       **end if**
19:     **end for**
20: **until** $done = 1$

---

---

**Algorithm 4** Modified Policy Iteration

---

1: **Input**: $\pi_0(s), \forall s \in \mathcal{S}$, m;
2: **Initialize**: $J(s) = 0, \pi(s) = \pi_0(s), \forall s \in \mathcal{S}$;
3: **repeat**
4:     **for** $k = 0, \cdots, m$ **do**
5:       **for** each $s \in \mathcal{S}$ **do**
6:         $J(s) = \sum_{s' \in \mathcal{S}} P_{ss'}(\pi(s))[r(s, \pi(s), s') + \gamma J(s')]$
7:       **end for**
8:     **end for**
9:     $done = 1$;
10:     **for** each $s \in \mathcal{S}$ **do**
11:       $b = \pi(s)$;
12:       $\pi(s) = \operatorname{argmax}_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} P_{ss'}(a)[r(s, a, s') + \gamma J(s')]$;
13:       **if** $b \neq \pi(s)$ **then**
14:         $done = 0$;
15:       **end if**
16:     **end for**
17: **until** $done = 1$

---