# Reinforcement Learning Assignment 1

Yifan Bai(260562421), Bhavya Patwa(260964036)

February 2020

## 1  Multi-Arm Bandits

We have K i.i.d distributions $\{R_i\}_{i=1}^K$ in $[0,1]$. Total arms pulls are $T$ so each arm is pulled $N_t(a) = \frac{T}{K}$ times for each action.

We need to guarantee that $u^* - u_{\hat{i}} \leq \epsilon, \forall \epsilon \geq 0$ with probability $1 - \delta, \delta \in (0,1)$.

We have the Hoeffding's inequality: $P[E[X] > \bar{X}_t + \mu] \leq e^{-2t\mu^2}$

Thus, applying Hoeffding's inequality to K-armed bandits, we get:

$$P[\mu^* - u_{\hat{i}} \leq \epsilon] = 1 - \delta$$
$$\therefore P[\mu^* - u_{\hat{i}} > \epsilon] = \delta$$
$$\therefore P[\mu^* > u_{\hat{i}}\epsilon] = \delta \leq e^{-2N_t(a)\epsilon^2}$$
$$\therefore \delta \leq e^{-2\frac{T}{K}\epsilon^2}$$
$$\therefore ln\delta \leq -2\frac{T}{K}\epsilon^2$$
$$\therefore -2T \geq ln(\delta)\frac{K}{\epsilon^2}$$
$$\therefore T \leq \frac{-Kln(\delta)}{2\epsilon^2}$$

Hence Proved.

## 2  Markov Decision Processes

### 2.1  Question i

$$V_M^\pi(s) = E[R_{t+1} + \gamma V^\pi(S_t + 1)|S_t = s]$$

Expanding, we get,

$$V_M^\pi(s)) = E[\overline{R}_{t+1} + \gamma\overline{R}_{t+2} + \gamma^2\overline{R}_{t+3} + ...|S_t = s]$$
$$= E[\overline{R}_{t+1} - \mathcal{N}(\mu, \sigma^2) + \gamma[\overline{R}_{t+2} - \mathcal{N}(\mu, \sigma^2)] + \gamma^2[\overline{R}_{t+2} - \mathcal{N}(\mu, \sigma^2)] + ...|S_t = s]$$

Now that the Gaussian is only of constant mean $\mu$ and variance $\sigma$, we can deduce,

$$= E[\overline{R}_{t+1} + \gamma\overline{R}_{t+2} + \gamma^2\overline{R}_{t+2} + ...|S_t = s] - E[-\mathcal{N}(\mu, \sigma^2) - \gamma\mathcal{N}(\mu, \sigma^2) - \gamma^2\mathcal{N}(\mu, \sigma^2) - ...|S_t = s]$$
$$= V_{\overline{(M)}}^\pi(s) - \mu[1 + \gamma + \gamma^2 + ...]$$

Since $\gamma \in [0,1)$, from summation of power series we have,

$$V_M^\pi(s) = V_{\overline{M}}^\pi(s) - \frac{\mu}{1 - \gamma}$$

### 2.2  Question ii

$$V_M^\pi(s)) = R^\pi + \gamma P_{\overline{M}}^\pi V_{\overline{M}}^\pi$$
$$= R^\pi + \gamma(\alpha P - \alpha Q + Q)V_{\overline{M}}^\pi$$
$$V_M^\pi(s)) = (1 - \gamma(\alpha P - \alpha Q + Q)^{-1}R^\pi$$

Then,

$$V_M^\pi(s)) - \gamma P_M^\pi V_M^\pi = V_{\overline{M}}^\pi(s)) - \gamma P_{\overline{M}}^\pi V_{\overline{M}}^\pi$$

$$= V_M^\pi(s)) - (\alpha\gamma P^\pi - \alpha\gamma Q^\pi + \gamma Q^\pi)V_M^\pi$$

Therefore,

$$V_M^\pi = \frac{1-\gamma P}{1-\alpha\gamma P - \alpha\gamma Q + \gamma Q}V_M^\pi$$

# 3  Policy Evaluation and Improvement

We have value function $\hat{V}$ such that $|V^*(s) - \hat{V}(s)| \le \epsilon$. Evaluating the greedy policy $\hat{V}_{\hat{V}}$ w.r.t $\hat{V}$ will give us:

$$\hat{V}_{\hat{V}}(s) = \arg\max[R_{t+1} + \gamma \sum_{s' \in S} P_{ss'}^a \hat{V}_{\hat{V}}(s')]$$

1. $0 < V^*(s) - \hat{V}(s) \le \epsilon$ Taking $L_{\hat{V}}(s) = V^*(s) - \hat{V}_{\hat{V}}(s)$, we get:

$$L_{\hat{V}}(s) = [R_{t+1} + \gamma \sum_{s' \in S} P_{ss'}^a V^*(s')] - \arg\max[R_{t+1} + \gamma \sum_{s' \in S} P_{ss'}^a \hat{V}_{\hat{V}}(s')]$$

$$= \gamma \sum_{s' \in S} P_{ss'}^a (V^*(s') - \hat{V}_{\hat{V}}(s'))]$$

The greedy policy $\hat{V}_{\hat{V}}(s) \ge \hat{V}(s), \forall s - (1)$

By assuming all $P$'s=1, this becomes a recursive form and it can be written as:

$$\le \gamma[\epsilon + \gamma(V^*(s'') - \hat{V}_{\hat{V}}(s''))]$$

(Max difference bound remains $\epsilon$) due to (1)

$$\le \gamma[\epsilon + \gamma\epsilon + \gamma^2\epsilon + ...]$$

This becomes a geometric sum and we get the result

$$\le \frac{\gamma\epsilon}{1-\gamma}$$

2. $0 > V^*(s) - \hat{V}(s) \ge -\epsilon$ Taking $L_{\hat{V}}(s) = V^*(s) - \hat{V}_{\hat{V}}(s)$, we get:

$$L_{\hat{V}}(s) = [R_{t+1} + \gamma \sum_{s' \in S} P_{ss'}^a V^*(s')] - \arg\max[R_{t+1} + \gamma \sum_{s' \in S} P_{ss'}^a \hat{V}_{\hat{V}}(s')]$$

$$= \gamma \sum_{s' \in S} P_{ss'}^a (V^*(s') - \hat{V}_{\hat{V}}(s'))]$$

The greedy policy $\hat{V}_{\hat{V}}(s) \ge \hat{V}(s), \forall s - (2)$

By assuming all $P$'s=1, this becomes a recursive form and it can be written as:

$$\ge \gamma[-\epsilon + \gamma(V^*(s'') - \hat{V}_{\hat{V}}(s''))]$$

(Max difference bound remains $-\epsilon$) due to (2)

$$\ge \gamma[-\epsilon - \gamma\epsilon - \gamma^2\epsilon + ...]$$

This becomes a geometric sum and we get the result

$$\ge \frac{-\gamma\epsilon}{1-\gamma}$$

So we get:

$$\frac{-\gamma\epsilon}{1-\gamma} \le V^*(s) - \hat{V}_{\hat{V}}(s) \le \frac{\gamma\epsilon}{1-\gamma}$$

adding $\frac{\gamma\epsilon}{1-\gamma}$ on

# 4  Statement of Contribution

Both team members contributed equally to the written part of this assignment, from brainstorming ideas, formulating solutions, to reviewing solutions and typing out the document. We hereby state that all the work presented in this report is that of the authors. The link to Google Colab notebook for programming question is: `https://colab.research.google.com/drive/13nmNXWho9nYPhSY7pQ4G8kvMPX1mwFQ5#scrollTo=uqgNsjF-piAl`