

# E1 277 Reinforcement Learning

## Assignment - 2

- KAWIN M

### Question ① (a)

Sample complexity of Q-value iteration:

$$\|Q_k - Q^*\|_\infty \leq \epsilon$$

$$\Rightarrow \|Q_k - T^*Q^*\|_\infty \leq \epsilon$$

[ $\because Q^*$  is a fixed point of  $T^*$

$$\Rightarrow Q^* = T^*Q^*]$$

$$\Rightarrow \|T^*Q_{k-1} - T^*Q^*\|_\infty \leq \epsilon$$

[Using Q-value iteration algorithm  $\Rightarrow Q_k = T^*Q_{k-1}$ ]

Expanding  $T^*Q_{k-1}$ ,  $T^*Q^*$  using the definitions of Bellman optimal operators.

$$\Rightarrow \left\| \gamma(s, a) + \gamma \int P(s' | s, a) \max_{a' \in A} Q_{k-1}(s', a') ds' \right.$$

$$\left. - \left[ \gamma(s, a) + \gamma \int P(s' | s, a) \max_{a' \in A} Q^*(s', a') ds' \right] \right\|_\infty \leq \epsilon$$

Cancelling out  $\gamma(s, a)$

$$\Rightarrow \left\| \gamma \int P(s' | s, a) \max_{a' \in A} Q_{k-1}(s', a') ds' - \right.$$

$$\left. \gamma \int P(s' | s, a) \max_{a' \in A} Q^*(s', a') ds' \right\|_\infty \leq \epsilon$$

$$\Rightarrow \left\| \gamma \left[ \int P(s' | s, a) \max_{a' \in A} Q_{k-1}(s', a') ds' - \right. \right.$$

$$\left. \int P(s' | s, a) \max_{a' \in A} Q^*(s', a') ds' \right\|_\infty \leq \epsilon$$

$$\Rightarrow \gamma \left\| \int P(s'|s, a) \max_{a' \in A} [Q_{k-1}(s', a') - Q^*(s', a')] ds' \right\|_{\infty} \leq \epsilon$$

$$\Rightarrow \gamma \left\| \max_{a' \in A} [Q_{k-1}(s', a') - Q^*(s', a')] \int P(s'|s, a) ds' \right\|_{\infty} \leq \epsilon$$

$$\Rightarrow \gamma \left\| \max_{a' \in A} [Q_{k-1}(s', a') - Q^*(s', a')] \cdot 1 \right\|_{\infty} \leq \epsilon$$

using definition of max-norm

$$\Rightarrow \gamma \max_{\substack{s' \in S, \\ a' \in A}} \left[ \max_{a' \in A} [Q_{k-1}(s', a') - Q^*(s', a')] \right] \leq \epsilon$$

$$\Rightarrow \gamma \max_{s' \in S, a' \in A} [Q_{k-1}(s', a') - Q^*(s', a')] \leq \epsilon$$

$$\Rightarrow \gamma \|Q_{k-1} - Q^*\|_{\infty} \leq \epsilon$$

$$\Rightarrow \gamma \|Q_{k-1} - Q^*\|_{\infty} \leq \epsilon$$

$$\Rightarrow \|Q_k - Q^*\|_{\infty} \leq \gamma \|Q_{k-1} - Q^*\|_{\infty} \leq \epsilon$$

Doing this iteratively

$$\leq \gamma^2 \|Q_{k-2} - Q^*\|_{\infty} \leq \epsilon$$

$\vdots$

$$\leq \gamma^k \|Q_0 - Q^*\|_{\infty} \leq \epsilon \rightarrow \textcircled{1}$$

We know that,

$$Q(s, a) = E[G_t | s_t = s, A_t = a] \leq G_t^{\max}$$

$$\text{and } G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$G_t^{\max} \triangleq R_{\max} + \gamma R_{\max} + \gamma^2 R_{\max} + \dots$$

$$= R_{\max} [1 + \gamma + \gamma^2 + \dots] \quad [R_{\max} = \text{Maximum possible reward value}]$$

$$= R_{\max} \cdot \frac{1}{1-\gamma}$$

$$G_t^{\max} = \frac{R_{\max}}{(1-\gamma)}$$

$$\therefore Q(s, a) \leq \frac{R_{\max}}{(1-\gamma)}$$

$$\Rightarrow \|Q_0 - Q^*\|_{\infty} \leq \frac{R_{\max}}{(1-\gamma)} \longrightarrow (2)$$

$$\Rightarrow \gamma^K \|Q_0 - Q^*\|_{\infty} \leq \frac{\gamma^K R_{\max}}{(1-\gamma)} \leq \epsilon$$

$$\Rightarrow \gamma^K \leq \frac{\epsilon (1-\gamma)}{R_{\max}}$$

Taking log on both sides

$$\log \gamma^K \leq \log \left( \frac{(1-\gamma)\epsilon}{R_{\max}} \right)$$

$$K \log \gamma \leq \log \left( \frac{(1-\gamma)\epsilon}{R_{\max}} \right)$$

$$K \leq \frac{\log \left( \frac{(1-\gamma)\epsilon}{R_{\max}} \right)}{\log \gamma}$$

$\therefore$  The upper bound on number of iterations  $K$  required to ensure that the error  $\|Q_K - Q^*\|_{\infty} \leq \epsilon$

is

$$K \leq \frac{\log\left(\frac{(1-\gamma)\epsilon}{R_{\max}}\right)}{\log \gamma}$$



### Question ①(b)

To compute the sample complexity of the policy iteration algorithm to ensure that  $\|Q^{\pi_k} - Q^*\|_{\infty} \leq \epsilon$

$$\|Q^{\pi_k} - Q^*\|_{\infty} \leq \epsilon$$

$$\Rightarrow \|Q^{\pi_k} - T^* Q^*\|_{\infty} \leq \epsilon \quad [\because Q^* \text{ is a fixed point of } T^* \\ \Rightarrow Q^* = T^* Q^*]$$

$$\Rightarrow \|T^{\pi_k} Q^{\pi_{k-1}} - T^* Q^*\|_{\infty} \leq \epsilon \quad [\text{Using Policy Evaluation} \\ \Rightarrow Q^{\pi_k} = T^{\pi_k} Q^{\pi_{k-1}}]$$

Expanding  $T^{\pi_k} Q^{\pi_{k-1}}$  and  $T^* Q^*$ ,

$$\Rightarrow \left\| r(s,a) + \gamma \int P(s'|s,a) \int \pi(a'|s') Q^{\pi_{k-1}}(s',a') da' ds' \right. \\ \left. - r(s,a) - \gamma \int P(s'|s,a) \max_{a' \in A} Q^*(s',a') ds' \right\|_{\infty} \leq \epsilon$$

Cancelling out  $r(s,a)$

$$\Rightarrow \left\| \gamma \int P(s'|s,a) \left[ \int \pi(a'|s') Q^{\pi_{k-1}}(s',a') da' \right. \right. \\ \left. \left. - \max_{a' \in A} Q^*(s',a') \right] ds' \right\|_{\infty} \leq \epsilon$$

$$\leq \gamma \left\| \int P(s'|s,a) \left[ \int \pi(a'|s') \max_{a' \in A} Q^{\pi_{k-1}}(s',a') da' \right. \right. \\ \left. \left. - \max_{a' \in A} Q^*(s',a') \right] ds' \right\|_{\infty} \leq \epsilon$$

$$\leq \gamma \left\| \int P(s'|s,a) \left[ \max_{a' \in A} Q^{\pi_{k-1}}(s',a') \int \pi(a'|s') da' \right. \right. \\ \left. \left. - \max_{a' \in A} Q^*(s',a') \right] ds' \right\|_{\infty} \leq \epsilon$$

$$\leq \gamma \left\| \int P(s' | s, a) \left[ \max_{a' \in A} (Q^{\pi_{k-1}}(s', a') - Q^*(s', a')) \right] ds' \right\|_{\infty} \leq \epsilon$$

$$\leq \gamma \left\| \max_{s', a'} (Q^{\pi_{k-1}}(s', a') - Q^*(s', a')) \int P(s' | s, a) ds' \right\|_{\infty} \leq \epsilon$$

$$\leq \gamma \left\| \max_{s', a'} (Q^{\pi_{k-1}}(s', a') - Q^*(s', a')) \cdot 1 \right\|_{\infty} \leq \epsilon$$

$$\leq \gamma \| Q^{\pi_{k-1}} - Q^* \|_{\infty} \leq \epsilon$$

$$\Rightarrow \| Q^{\pi_k} - Q^* \|_{\infty} \leq \gamma \| Q^{\pi_{k-1}} - Q^* \|_{\infty} \leq \epsilon$$

Doing this iteratively, we get

$$\leq \gamma^2 \| Q^{\pi_{k-2}} - Q^* \|_{\infty} \leq \epsilon$$

$$\vdots$$

$$\leq \gamma^K \| Q^{\pi_0} - Q^* \|_{\infty} \leq \epsilon$$

using Equation (2) from Question 1(a)

$$\Rightarrow \| Q^{\pi_k} - Q^* \|_{\infty} \leq \gamma^K \frac{R_{\max}}{(1-\gamma)} \leq \epsilon$$

$$\Rightarrow \gamma^K \leq \frac{\epsilon(1-\gamma)}{R_{\max}}$$

Taking log on both sides

$$\Rightarrow K \log \gamma \leq \log \left( \frac{\epsilon(1-\gamma)}{R_{\max}} \right)$$

$$\Rightarrow K \leq \frac{\log \left( \frac{\epsilon(1-\gamma)}{R_{\max}} \right)}{\log \gamma}$$

### Question (1)(c)

Considering all the non-stationary policies  $\tilde{\pi}_1, \tilde{\pi}_2, \tilde{\pi}_3, \dots$  over the state space  $S$  and action space  $A$ .

Let the order of the value functions corresponding to the non-stationary policies be

$$V^{\tilde{\pi}_1} \leq V^{\tilde{\pi}_2} \leq V^{\tilde{\pi}_3} \leq \dots$$

It is given that

$$V^+ = \sup_{\tilde{\pi}} V^{\tilde{\pi}}$$

$$\therefore V^{\tilde{\pi}_1} \leq V^{\tilde{\pi}_2} \leq V^{\tilde{\pi}_3} \leq \dots \leq V^+ \longrightarrow (1)$$

But, we also know that

$$V^{\tilde{\pi}_1} \leq V^{\tilde{\pi}_2} \leq V^{\tilde{\pi}_3} \leq \dots \leq V^* \longrightarrow (2)$$

From (1) & (2)

$$V^+ \leq V^* \quad \text{and} \quad V^* \leq V^+$$

$$\Rightarrow V^+ = V^*$$

Hence Proved

Conclusion:

It can be concluded that even in the continuous state and action space and with non-stationary policies, repeated application of the Bellman operator with supremum value converges to the optimal value.



Question (1)(d)

Given:  $(\hat{T}^\pi v)(s) = \gamma^\pi(s) \int P(s'|s) v(s') ds'$

i) Monotonicity:

To Prove: For any  $v, \bar{v} \in \mathbb{R}^n$ , such that, let

$$v(i) \leq \bar{v}(i) \quad \forall 1 \leq i \leq n$$

then,  $(\hat{T}_K^\pi v)(i) \leq (\hat{T}_K^\pi \bar{v})(i)$

Proof using Induction:

Base case:  $K=1$

$$\begin{aligned} (\hat{T}^\pi v)(i) &= \gamma^\pi(i) \int P(s'|i) v(s') ds' \\ &\leq \gamma^\pi(i) \int P(s'|i) \bar{v}(s') ds' \quad [\because v(i) \leq \bar{v}(i)] \\ &\leq (\hat{T}^\pi \bar{v})(i) \end{aligned}$$

Hence Proved

Induction Hypothesis: Let us assume that for any arbitrary  $K$ ,

$$(\hat{T}_K^\pi v)(i) \leq (\hat{T}_K^\pi \bar{v})(i) \text{ is true}$$

To Prove:  $(\hat{T}_{K+1}^\pi v)(i) \leq (\hat{T}_{K+1}^\pi \bar{v})(i)$

$$\begin{aligned} \text{Proof: } (\hat{T}_{K+1}^\pi v)(i) &= (\hat{T}^\pi v_K^\pi)(i) \\ &= \gamma^\pi(i) \int P(s'|i) v_K^\pi(s') ds' \\ &= \gamma^\pi(i) \int P(s'|i) (\hat{T}_K^\pi v)(s') ds' \end{aligned}$$



$$\leq \gamma^\pi(i) \int P(s'|i) (\hat{T}_k^\pi \bar{V})(s') ds' \quad [\because \text{Induction Hypothesis}]$$

$$\leq \gamma^\pi(i) \int P(s'|i) (\bar{V}_k^\pi)(s') ds'$$

$$\leq (\hat{T}_{k+1}^\pi \bar{V})(i)$$

Hence proved that the given Bellman operator is monotonic using Induction

(ii) contraction:

$$\text{To Prove: } \|\hat{T}^\pi V - \hat{T}^\pi \bar{V}\|_\infty \leq \beta \|V - \bar{V}\|_\infty,$$

for any  $V, \bar{V} \in \mathbb{R}^n$  and a scalar  $0 < \beta < 1$

Proof:

$$\begin{aligned} \hat{T}^\pi V - \hat{T}^\pi \bar{V} &= \gamma^\pi(s) \int P(s'|s) V(s') ds' \\ &\quad - \gamma^\pi(s) \int P(s'|s) \bar{V}(s') ds' \end{aligned}$$

$$= \gamma^\pi(s) \int P(s'|s) [V(s') - \bar{V}(s')] ds'$$

$$|\hat{T}^\pi V - \hat{T}^\pi \bar{V}| = |\gamma^\pi(s) \int P(s'|s) [V(s') - \bar{V}(s')] ds'|$$

$$= \gamma^\pi(s) \int P(s'|s) |V(s') - \bar{V}(s')| ds'$$

Taking max on both sides  
 $i=1..n$

$$\max_{i=1..n} |\hat{T}^\pi V - \hat{T}^\pi \bar{V}| \leq \gamma^\pi(s) \int P(s'|s) \max_{s=1..n} |V(s') - \bar{V}(s')| ds'$$

$$\Rightarrow \|\hat{T}^\pi V - \hat{T}^\pi \bar{V}\|_\infty \leq \gamma^\pi(s) \max_{s=1..n} |V(s') - \bar{V}(s')| \int P(s'|s) ds'$$

$$\|\hat{T}^\pi v - \hat{T}^\pi \bar{v}\|_\infty \leq \gamma^\pi(s) \|v - \bar{v}\|_\infty \cdot 1$$

$$\Rightarrow \|\hat{T}^\pi v - \hat{T}^\pi \bar{v}\|_\infty \leq \beta \|v - \bar{v}\|_\infty$$

$$\text{where } \beta = \gamma^\pi(s)$$

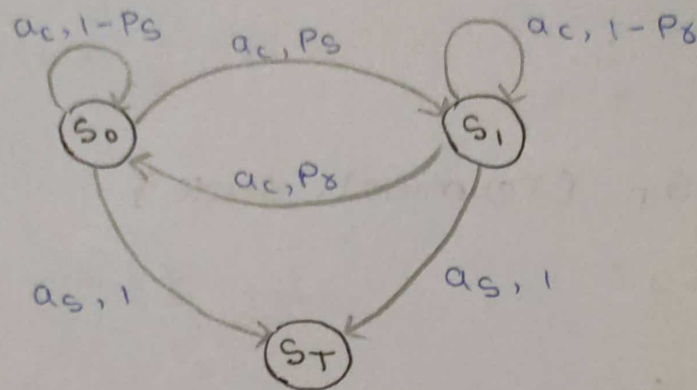
$\therefore$  The given operator is a contraction  
given  $0 < \gamma^\pi(s) < 1$

$$\text{Contraction factor} = \gamma^\pi(s)$$

Hence Proved

Question (2) (a)

Transition Probabilities for the given MDP



$$P(s_0 | s_0, a_s) = 0$$

$$P(s_1 | s_0, a_s) = 0$$

$$P(s_T | s_0, a_s) = 1$$

$$P(s_0 | s_0, a_c) = 1 - P_s$$

$$P(s_1 | s_0, a_c) = P_s$$

$$P(s_T | s_0, a_c) = 0$$

$$P(s_0 | s_1, a_s) = 0$$

$$P(s_1 | s_1, a_s) = 0$$

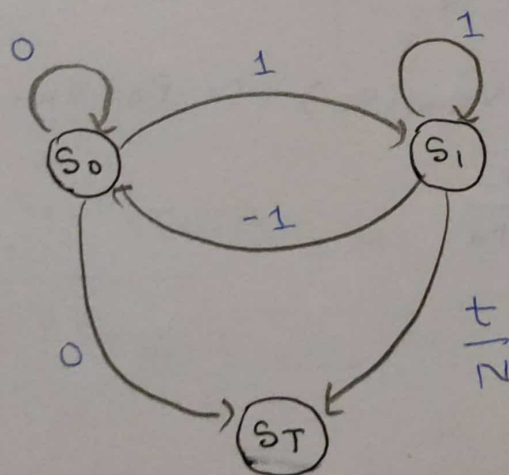
$$P(s_T | s_1, a_s) = 1$$

$$P(s_0 | s_1, a_c) = P_s$$

$$P(s_1 | s_1, a_c) = 1 - P_s$$

$$P(s_T | s_1, a_c) = 0$$

Question (2) (b)



$$t = 1, 2, \dots, N$$



Question (2) (c)

By Dynamic Programming,

$$V_k(s) = \max_a \sum_{s', r} P(s', r | s, a) [r + V_{k+1}(s')]$$

For state  $s_T$  [Terminal state]

$$V_k(s_T) = 0, \quad \forall k = 1, \dots, N$$

For state  $s_0$

$$V_k(s_0) = \max \begin{cases} P_s [1 + V_{k+1}(s_1)] + (1 - P_s) [0 + V_{k+1}(s_0)], \\ 1 [0 + V_{k+1}(s_T)] \end{cases}$$

$$= \max \{ P_s + P_s V_{k+1}(s_1) + (1 - P_s) V_{k+1}(s_0), 0 \}$$

For state  $s_1$

$$V_k(s_1) = \max \begin{cases} P_\delta [-1 + V_{k+1}(s_0)] + (1 - P_\delta) [1 + V_{k+1}(s_1)], \\ 1 \left[ \frac{k}{N} + V_{k+1}(s_T) \right] \end{cases}$$

$$= \max \begin{cases} P_\delta V_{k+1}(s_0) + (1 - P_\delta) V_{k+1}(s_1) + (1 - 2P_\delta), \\ \frac{k}{N} \end{cases}$$

Question (3)(a)

To show: Maximizing the expected total sum of rewards  $\mathbb{E}\left[\sum_{t=1}^T x_{t, I_t}\right]$  is equivalent to minimizing the regret.

$$\Rightarrow \max \mathbb{E}\left[\sum_{t=1}^T x_{t, I_t}\right] = \min \left( -\mathbb{E}\left[\sum_{t=1}^T x_{t, I_t}\right] \right)$$

$$= \min \left[ T\mu^* - \mathbb{E}\left[\sum_{t=1}^T x_{t, I_t}\right] \right]$$

$$= \min \left[ T\mu^* - \mathbb{E}\left[\sum_{i=1}^K \sum_{t=1}^T x_{ti} \mathbb{1}_{\{I_t=i\}}\right] \right]$$

Dividing and Multiplying by  $\sum_{t=1}^T \mathbb{1}_{\{I_t=i\}}$

$$= \min \left( T\mu^* - \mathbb{E}\left[ \sum_{i=1}^K \left[ \frac{\sum_{t=1}^T x_{ti} \mathbb{1}_{\{I_t=i\}}}{\sum_{t=1}^T \mathbb{1}_{\{I_t=i\}}} \cdot \sum_{t=1}^T \mathbb{1}_{\{I_t=i\}} \right] \right] \right)$$

$$\text{Here, } \frac{\sum_{t=1}^T x_{ti} \mathbb{1}_{\{I_t=i\}}}{\sum_{t=1}^T \mathbb{1}_{\{I_t=i\}}} = \mu_i \quad \text{and} \quad \sum_{t=1}^T \mathbb{1}_{\{I_t=i\}} = N_i(T)$$

$$\therefore = \min \left( T\mu^* - \mathbb{E}\left[\sum_{i=1}^K \mu_i \cdot N_i(T)\right] \right)$$

$$= \min (\text{Regret}(T))$$

Hence shown that  $\max \mathbb{E}\left[\sum_{t=1}^T x_{t, I_t}\right] = \min \text{Regret}(T)$

Question (3)(b)(i)

Given:  $K = 2$

Rewards are bounded in  $[0, 1]$

Arm 1 is optimal  $\Rightarrow \mu^* = \mu_1$

$$\Delta = \mu_1 - \mu_2$$

To prove:  $\text{Regret}(T) \leq m\Delta + \Delta(T - 2m) \mathbb{E} \left[ \mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)} \right]$

Proof:

$$\text{Regret}(T) = T\mu^* - \mathbb{E} \left[ \sum_{i=1}^K \mu_i N_i(T) \right]$$

Phase 1:

In phase 1 for time  $t = 1, 2, \dots, 2m$  the 2 arms are played in a round-robin fashion; i.e. each arm is played  $m$  times ( $\because \frac{T}{2} = \frac{2m}{2} = m$ )

$$\begin{aligned} \therefore \text{Regret}(T_1) &= (2m)\mu_1 - \mathbb{E} \left[ \sum_{i=1}^2 \mu_i N_i(T) \right] \\ &= (2m)\mu_1 - \mathbb{E} [\mu_1 N_1(T) + \mu_2 N_2(T)] \end{aligned}$$

Here,  $N_1(T) = N_2(T) = m$

$$\begin{aligned} &= 2m\mu_1 - \mathbb{E} [\mu_1 \cdot m + \mu_2 m] \\ &= 2m\mu_1 - \mu_1 m - \mu_2 m \\ &= m\mu_1 - m\mu_2 \\ &= m(\mu_1 - \mu_2) \end{aligned}$$

$$\text{Regret}(T_1) = m\Delta \longrightarrow \textcircled{1}$$

$\therefore$  Regret in phase 1 is  $m\Delta$



Phase 2:

In phase 2 for time  $t = 2m+1, 2m+2, \dots, T$ , the arm with the best sample mean till  $t = 2m$  is played.

$$\therefore \text{Total time } T_2 = T - (2m+1) - 1 \\ = T - 2m$$

Number of times arm 1 is played:

$$N_1(T_2) = T_2 * \mathbb{1}_{\hat{\mu}_1(2m) > \hat{\mu}_2(2m)} \\ = (T-2m) \mathbb{1}_{\hat{\mu}_1(2m) > \hat{\mu}_2(2m)}$$

Number of times arm 2 is played:

$$N_2(T_2) = T_2 * \mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)} \\ = (T-2m) \mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)}$$

$$\therefore \text{Regret}(T_2) = (T-2m)\mu_1 - \mathbb{E} \left[ \mu_1(T-2m) \mathbb{1}_{\hat{\mu}_1(2m) > \hat{\mu}_2(2m)} \right. \\ \left. + \mu_2(T-2m) \mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)} \right]$$

$$= (T-2m)\mu_1 - \mathbb{E} \left[ \mu_1(T-2m) (1 - \mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)}) \right. \\ \left. + \mu_2(T-2m) \mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)} \right]$$

$$= (T-2m)\mu_1 - \mathbb{E} \left[ \mu_1(T-2m) - \mu_1(T-2m) \mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)} \right. \\ \left. + \mu_2(T-2m) \mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)} \right]$$

$$= (T-2m)\mu_1 - \mu_1(T-2m) - \mathbb{E} \left[ (\mu_2 - \mu_1) \mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)} \right] \\ \cdot (T-2m)$$

cancelling out  $(T-2m)\mu_1$

$$\leq -\mathbb{E}[-\Delta(T-2m) \mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)}]$$

$$\leq \Delta(T-2m) \mathbb{E}[\mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)}]$$

$\therefore$  Regret in Phase 2 is

$$\text{Regret}(T_2) \leq \Delta(T-2m) \mathbb{E}[\mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)}] \rightarrow (2)$$

From (1) & (2)

Total Regret

$$\text{Regret}(T) \leq \text{Regret}(T_1) + \text{Regret}(T_2)$$

$$\leq m\Delta + \Delta(T-2m) \mathbb{E}[\mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)}]$$

Hence Proved

Question (3) (b) (ii)

$$\text{Regret}(T) \leq m\Delta + \Delta(T-2m) \mathbb{E}[\mathbb{1}_{\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)}]$$

We know that,

$$\mathbb{E}[\mathbb{1}_A] = P(A)$$

$$\therefore \text{Regret}(T) \leq m\Delta + \Delta(T-2m)P(\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m)) \rightarrow (3)$$

Here, considering  $P(\hat{\mu}_1(2m) \leq \hat{\mu}_2(2m))$

Expanding  $\hat{\mu}_1(2m)$  and  $\hat{\mu}_2(2m)$

$$P\left(\frac{\sum_{s=1}^{2m} X_{S1} \mathbb{1}_{\{I_S=1\}}}{\sum_{s=1}^{2m} \mathbb{1}_{\{I_S=1\}}} \leq \frac{\sum_{s=1}^{2m} X_{S2} \mathbb{1}_{\{I_S=2\}}}{\sum_{s=1}^{2m} \mathbb{1}_{\{I_S=2\}}}\right)$$

$$\text{Here, } \sum_{s=1}^{2m} \mathbb{1}_{\{I_S=1\}} = N_1(2m) = m$$

$$\sum_{s=1}^{2m} \mathbb{1}_{\{I_S=2\}} = N_2(2m) = m \quad (\text{Phase 1})$$

$$\Rightarrow P\left(\frac{\sum_{s=1}^{2m} X_{S1} \mathbb{1}_{\{I_S=1\}}}{m} \leq \frac{\sum_{s=1}^{2m} X_{S2} \mathbb{1}_{\{I_S=2\}}}{m}\right)$$

Adding and Subtracting  $\frac{\mu_1}{m}$  on LHS and

$\frac{\mu_2}{m}$  on RHS

$$\Rightarrow P\left(\frac{\sum_{s=1}^{2m} X_{S1} \mathbb{1}_{\{I_S=1\}} - \mu_1}{m} \leq \frac{\sum_{s=1}^{2m} X_{S2} \mathbb{1}_{\{I_S=2\}} - \mu_2}{m} + \frac{\mu_2 - \mu_1}{m}\right)$$



Rearranging the terms

$$P\left(\frac{\mu_1 - \mu_2}{m} \leq \frac{\left[\sum_{s=1}^{2m} x_{s,2} \mathbb{1}_{\{I_s=2\}} - \mu_2\right] - \left[\sum_{s=1}^{2m} x_{s,1} \mathbb{1}_{\{I_s=1\}} - \mu_1\right]}{m}\right)$$

w.k.t  $E[x_1] = \mu_1$  and  $E[x_2] = \mu_2$

$$P\left(\frac{\Delta}{m} \leq \frac{1}{m} \left( \left( \sum_{s=1}^{2m} x_{s,2} \mathbb{1}_{\{I_s=2\}} - \mu_2 \right) + \left( \sum_{s=1}^{2m} x_{s,1} \mathbb{1}_{\{I_s=1\}} - \mu_1 \right) \right) \right)$$

$$\Rightarrow P\left(\frac{\Delta}{m} \leq \frac{1}{m} \left( \left( \sum_{s=1}^{2m} x_{s,2} \mathbb{1}_{\{I_s=2\}} - E[x_2] \right) + \left( \sum_{s=1}^{2m} x_{s,1} \mathbb{1}_{\{I_s=1\}} - E[x_1] \right) \right) \right)$$

$$\Rightarrow P\left(\frac{\Delta}{m} \leq \frac{1}{m} \left( \sum_{s=1}^{2m} x_{s,I_s} - E[x_{I_s}] \right) \right)$$

Dividing by 2 on both sides

$$\Rightarrow P\left(\frac{\Delta}{2m} \leq \frac{1}{2m} \left( \sum_{s=1}^{2m} x_{s,I_s} - E[x_{I_s}] \right) \right)$$

Here, as  $\{x_{s,I_s}\}_s$  is i.i.d and  $x_{s,I_s} \in [0, 1]$

$\therefore \Rightarrow a=0$ ,  $b=1$  and

$$\epsilon = \frac{\Delta}{2m}, \quad n = 2m$$

$\therefore$  Using Chernoff-Hoeffding's bound,

$$\begin{aligned} \Rightarrow P\left(\frac{\Delta}{2m} \leq \frac{1}{2m} \left( \sum_{s=1}^{2m} x_{s,I_s} - E[x_{I_s}] \right) \right) &\leq e^{-2 \cdot 2m \cdot \frac{\Delta^2}{4m^2} / (1-0)^2} \\ &\leq e^{-\left(\frac{1}{m}\right)m\Delta^2} \end{aligned}$$

$$\leq e^{-cm\Delta^2}, \text{ where } c > 0$$

Substituting in (3)

$$\begin{aligned} \text{Regret}(T) &\leq m\Delta + \Delta(T-2m)e^{-cm\Delta^2} \\ &\leq m\Delta + \Delta Te^{-cm\Delta^2} \end{aligned}$$

Hence shown, that  $\text{Regret}(T) \leq m\Delta + \Delta Te^{-cm\Delta^2}$

$$\text{Let } m = \frac{\log T}{c\Delta^2}$$

$$\text{then, } \text{Regret}(T) \leq \frac{\log T \cdot \Delta}{c\Delta^2} + \Delta(T-2m)e^{-\frac{c \cdot \log T}{c\Delta^2} \cdot \Delta}$$

$$\leq \frac{\log T}{c\Delta} + \Delta Te^{-\log T}$$

$$\leq \frac{\log T}{c\Delta} + \Delta T e^{\log T^{-1}}$$

$$\leq \frac{\log T}{c\Delta} + \Delta T \cdot T^{-1}$$

$$\leq \frac{\log T}{c\Delta} + \Delta$$

$$\text{Regret}(T) = O\left(\frac{1}{\Delta} \log T\right)$$

Question (3) (b) (iii)

Difficulty in implementing the algorithm to obtain a sub-linear regret:

i) It takes linear time on  $T$  to compute  $\mu^*$ .

ii) It is difficult to compute the expected reward obtained by the algorithm.