

DDA 4230 Assignment 1

Name: Xiang Fei ; ID: 120090414

P1.

$$1. \pi_1 : s_0 \rightarrow a_1, s_1 \rightarrow a_0, s_2 \rightarrow a_0, s_3 \rightarrow a_0$$

$$\pi_2 : s_0 \rightarrow a_2, s_1 \rightarrow a_0, s_2 \rightarrow a_0, s_3 \rightarrow a_0$$

$$2. V^*(s_0) = \max_{a \in \{a_1, a_2\}} 0 + r \sum_{s'} P(s'|s, a) V^*(s') = r \max\{V^*(s_1), V^*(s_2)\}$$

$$V^*(s_1) = \max_{a \in \{a_0\}} 0 + r \sum_{s'} P(s'|s, a) V^*(s') = r[(1-p)V^*(s_1) + pV^*(s_3)]$$

$$V^*(s_2) = \max_{a \in \{a_0\}} 1 + r \sum_{s'} P(s'|s, a) V^*(s') = 1 + r[(1-q)V^*(s_0) + qV^*(s_3)]$$

$$V^*(s_3) = \max_{a \in \{a_0\}} 10 + r \sum_{s'} P(s'|s, a) V^*(s') = 10 + rV^*(s_0)$$

3. No.

$$\text{if } r=0, \text{ since } \pi^*(s_0) = \arg\max_{a \in \{a_1, a_2\}} r \sum_{s'} P(s'|s, a) V^*(s')$$

so π^* is not unique, then $\pi^*(s_0)=a_2$ can not be sure.

4. No.

$$V^*(s_2) = 1 + r[(1-q)V^*(s_0) + qV^*(s_3)] \geq 1$$

$$\Rightarrow \text{iff } V^*(s_1) > V^*(s_2) \geq 1, \forall r \in [0, 1) \text{ and } p > 0 \rightarrow \pi^*(s_0) = a_1$$

$$\Rightarrow V^*(s_1) = r[(1-p)V^*(s_1) + pV^*(s_3)] = \frac{r p V^*(s_3)}{1-r(1-p)}$$

$$\Rightarrow \exists r \text{ such that } V^*(s_1) < 1.$$

P2.

$$\begin{aligned}
 1. \quad \|BV - BV'\|_\infty &= \left\| \max_a \left[R(s, a) + \gamma \sum_{s_j \in S} P(s_j | s, a) V(s_j) \right] - \max_a \left[R(s, a) + \gamma \sum_{s_j \in S} P(s_j | s, a) V'(s_j) \right] \right\|_\infty \\
 &\leq \max_a \left\| \left[R(s, a) + \gamma \sum_{s_j \in S} P(s_j | s, a) V(s_j) \right] - \left[R(s, a) + \gamma \sum_{s_j \in S} P(s_j | s, a) V'(s_j) \right] \right\|_\infty \\
 &\leq \gamma \max_a \left\| \sum_{s_j \in S} P(s_j | s, a) V(s_j) - \sum_{s_j \in S} P(s_j | s, a) V'(s_j) \right\|_\infty \\
 &= \gamma \max_a \left\| \left[\sum_{s_j \in S} P(s_j | s, a) (V(s_j) - V'(s_j)) \right] \right\|_\infty \\
 &\leq \gamma \max_{a, s_i} \sum_{s_j \in S} P(s_j | s_i, a) |V(s_j) - V'(s_j)| \\
 &\leq \gamma \max_{a, s_i} \sum_{s_j \in S} P(s_j | s_i, a) |V - V'| = \gamma \|V - V'\|_\infty \quad , \text{ Q.E.D.}
 \end{aligned}$$

2.1, basic case: $\|V_2 - V_1\|_\infty = \|BV_1 - BV_0\|_\infty \leq \gamma \|V_1 - V_0\|_\infty$

inductive assumption: $\|V_{k+1} - V_k\|_\infty \leq \gamma^k \|V_1 - V_0\|_\infty$

inductive step: $\|V_{k+2} - V_{k+1}\|_\infty = \|BV_{k+1} - BV_k\|_\infty \leq \gamma \|V_{k+1} - V_k\|_\infty \leq \gamma^{k+1} \|V_1 - V_0\|_\infty$

$\Rightarrow \|V_{n-1} - V_n\|_\infty \leq \gamma^n \|V_1 - V_0\|_\infty \quad , \text{ Q.E.D.}$

2.2, $\|V_{n+c} - V_n\| \leq \|V_{n+c} - V_{n+c-1}\| + \|V_{n+c-1} - V_n\|$

$$\begin{aligned}
 &\leq \|V_{n+c} - V_{n+c-1}\| + \dots + \|V_{n+1} - V_n\| \\
 &\leq \gamma^{n+c-1} \|V_1 - V_0\| + \gamma^{n+c-2} \|V_1 - V_0\| + \dots + \gamma^n \|V_1 - V_0\| \\
 &= \frac{\gamma^n (1 - \gamma^c)}{1 - \gamma} \|V_1 - V_0\| \leq \frac{\gamma^n}{1 - \gamma} \|V_1 - V_0\| \quad , \text{ Q.E.D.}
 \end{aligned}$$

3, let $k < n$, $m = n + c$, $\varepsilon = \frac{\gamma^k}{1 - \gamma} \|V_1 - V_0\|_\infty$

$$\begin{aligned}
 \|V_m - V_n\|_\infty &\leq \frac{\gamma^n}{1 - \gamma} \|V_1 - V_0\|_\infty \\
 &< \frac{\gamma^k}{1 - \gamma} \|V_1 - V_0\|_\infty = \varepsilon
 \end{aligned}$$

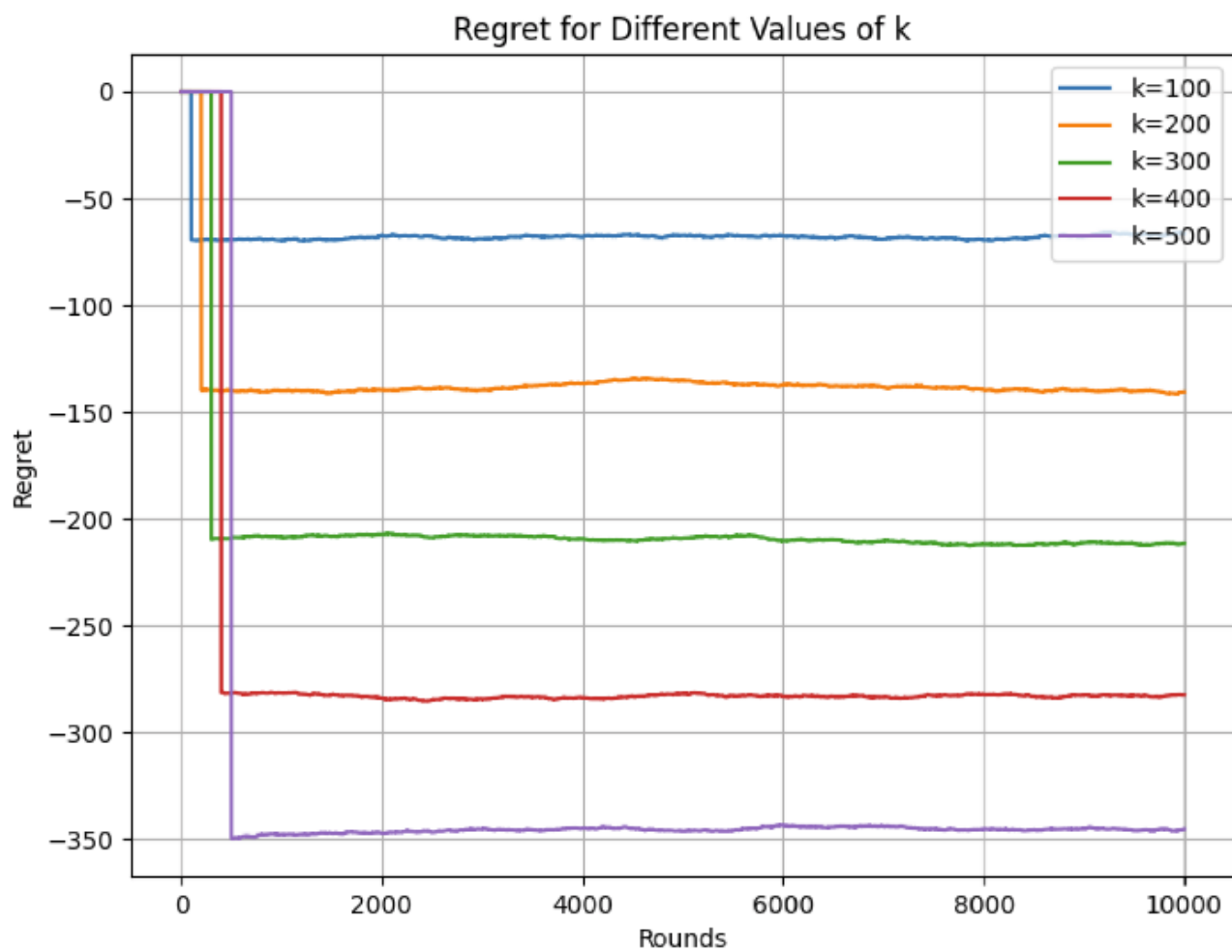
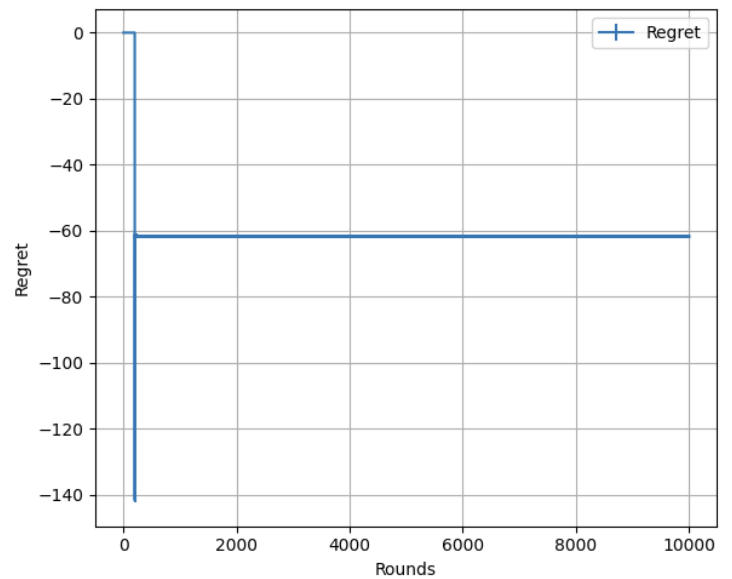
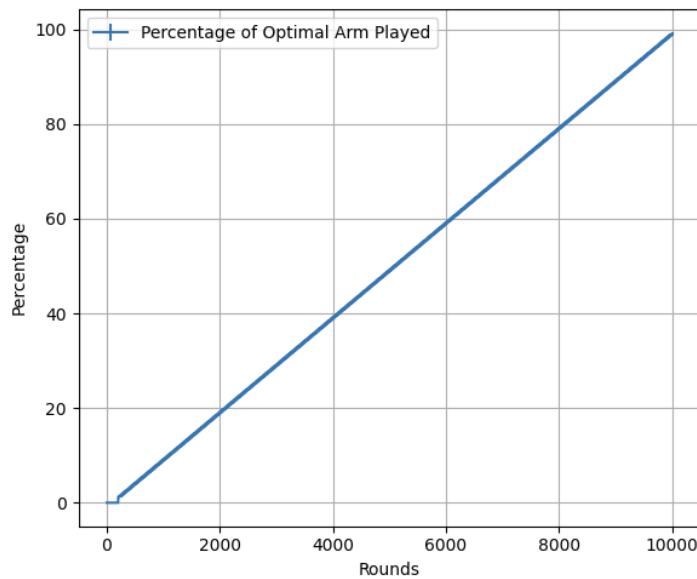
4. Assume two fixed points V and V'

$$\|V - V'\|_\infty = \|BV - BV'\|_\infty \leq \gamma \|V - V'\|_\infty$$

since $0 < \gamma < 1 \Rightarrow \|V - V'\|_\infty = 0 \Rightarrow V = V'$

therefore, only one fixed point

P3



1. ETC with Optimal k :

- In the case of ETC with an optimal choice of k , we would ideally choose k such that it minimizes the gap-dependent regret. However, in practice, we often lack knowledge of the underlying gap, so this choice is challenging.
- Theoretical findings suggest that for the ETC algorithm with an optimal k , the regret should be minimized because we initially explore both arms for a limited number of rounds (k) and then commit to the arm that appears to be the best based on the collected data.
- In the implementation, we experimented with different values of k , and if we were able to choose an optimal k that closely matched the true underlying gap, we would expect to see low regret. The regret would decrease quickly in the commitment phase as we concentrate on the better arm.

2. ETC with Heuristic Choice for k :

- In practice, we often have to use heuristics or rules of thumb to select the value of k because we typically do not know the true gap.
- When we choose k heuristically, it may not necessarily be optimal. The performance of ETC with a suboptimal k can vary depending on how different the chosen k is from the optimal one.
- If the chosen k is too small, the exploration phase might be insufficient, leading to suboptimal arm selection in the commitment phase, resulting in higher regret. Conversely, if k is too large, it could lead to prolonged exploration and slower convergence to the optimal arm.
- In the implementation, we experimented with different values of k , and we can observe that as k deviates from the theoretically optimal value, the regret might increase. For example, if k is far from the optimal value, the algorithm may spend too much time in exploration or may commit to a suboptimal arm for too long, leading to higher regret.

Overall, the choice of k in the ETC algorithm is critical, and the performance depends on how well it aligns with the true gap between the arms. In practice, it's often challenging to determine the optimal value of k without prior knowledge of the gap, so heuristics are used. The results obtained through experimentation with different values of k help in understanding the trade-off between exploration and commitment and provide insights into the algorithm's performance in the absence of perfect information about the gap.