

CS 747

Programming Assignment 1

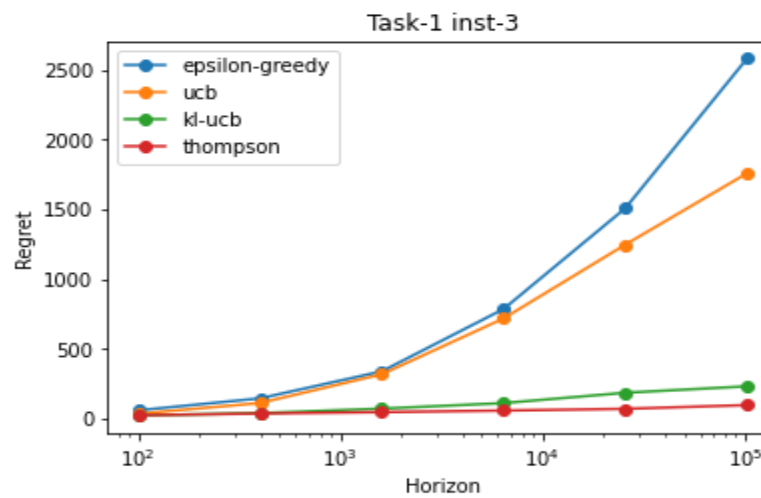
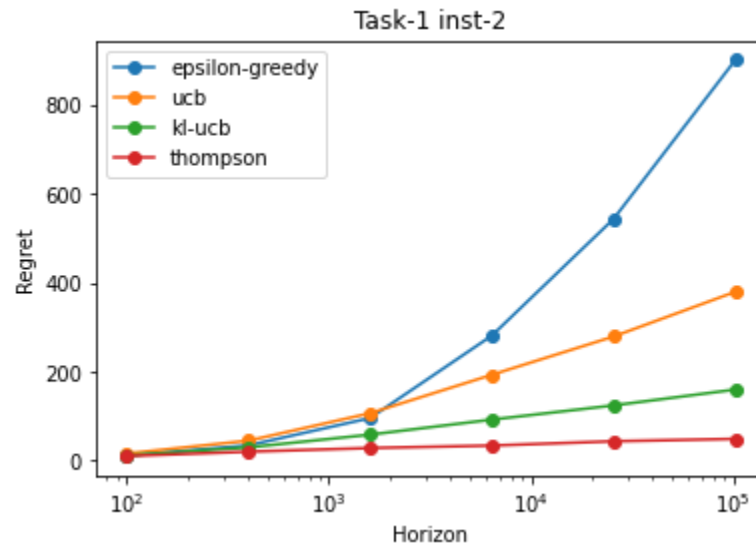
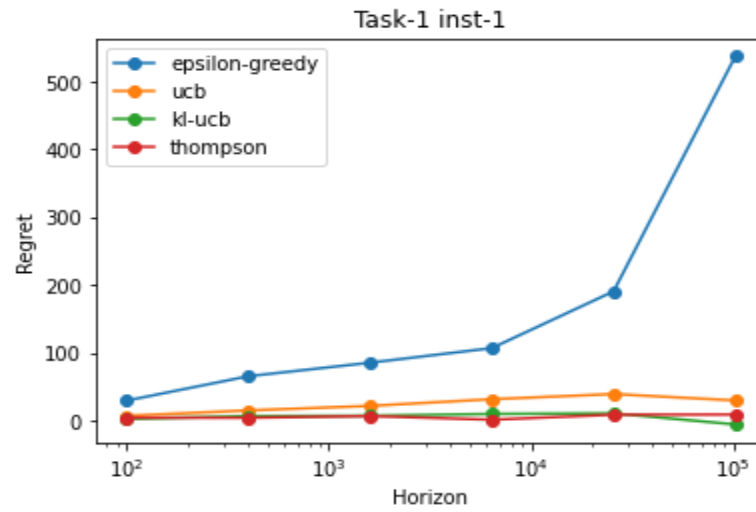
Rushikesh J. Metkar

19D070034

CONTENTS :

- **TASK1**
- **TASK2**
- **TASK3**
- **TASK4**

TASK 1



I have implemented the following sampling algorithms:

(1) epsilon-greedy, (2) UCB, (3) KL-UCB, and (4) Thompson Sampling

I've assumed that each arm is pulled at least once.

In Epsilon-greedy algorithm, I'm sampling an arm uniformly at random; with probability epsilon and sampling an arm with the highest empirical mean with probability $1 - \epsilon$

I know the most optimal arm and by using its true mean (p_{star}) I've computed regret as $\text{REG} = [hz \cdot p_{\text{star}} - \text{REW (cumulative reward)}]$

The Epsilon-greedy algorithm gives the highest regret for all instances.

In the UCB algorithm, I'm sampling an arm which has the highest ucb value.

$\text{ucb_arm} = \text{empirical_mean_arm} + \text{exploration_bonus}$ where

exploration_bonus is defined as the square root of $2 \times \ln(\text{total_pulls}) / \text{pulls_arm}$

REG is computed similarly to the first case.

The UCB algorithm gives the second highest regret for all instances.

In the KL-UCB algorithm, I'm sampling an arm which has the highest kl-ucb value.

$\text{ucb_kl}(a) = \max\{q \in [p^*(a), 1] \text{ such that } \text{KL}(p^*(a), q) \leq \ln(t) + c \ln(\ln(t))\}$, where $c \geq 3$

I've made two functions KL and solve_q where the first function finds the KL divergence between two inputs and the second function gives the value of q that follows KL inequality.

REG is computed similarly to the first case.

Regret is very low as compared to the first two cases.

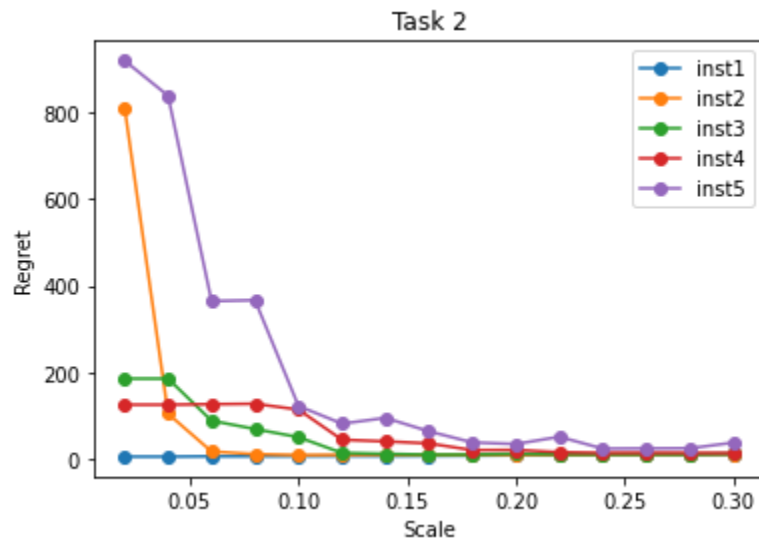
In the Thompson sampling algorithm, for every arm a, I'm drawing a sample $x_t(a) \sim \text{Beta}(s_t(a) + 1, f_t(a) + 1)$ where $s_t(a)$ are successes and $f_t(a)$ are failures and choosing the arm with maximal $x_t(a)$.

In this algorithm I've pulled each arm for 1 time in the start and then selected an arm depending on the Beta distribution values.

REG is computed similarly to the first case.

Regret is very low as compared to the first two cases.

TASK 2



The value of c that gave the lowest regret for each of the five instances is **0.26**

For each instance as we increase the scale, regret value decreases.

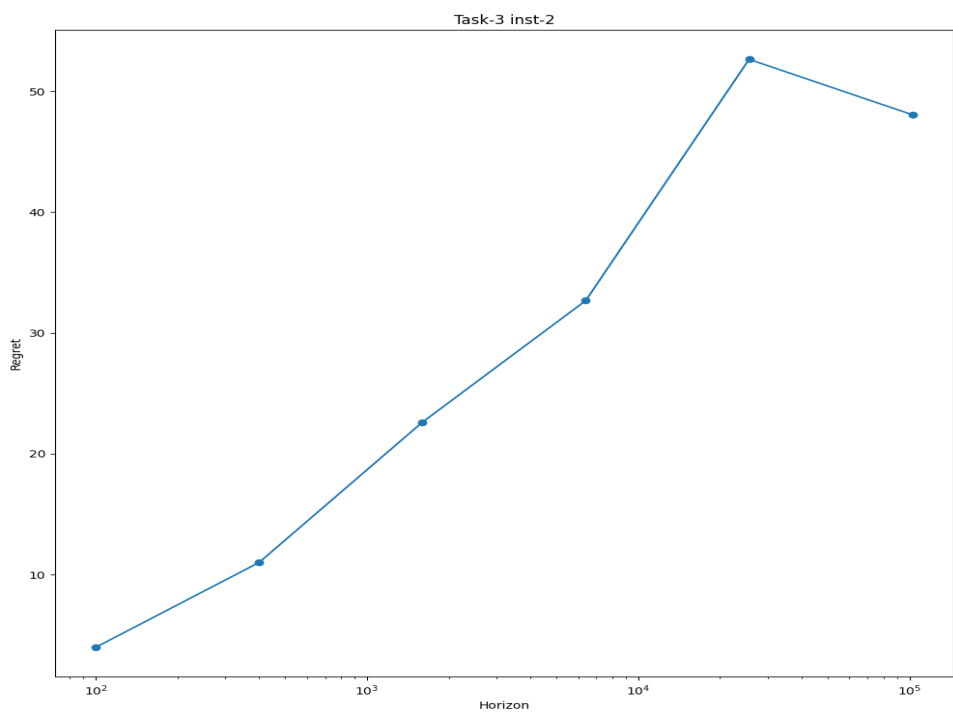
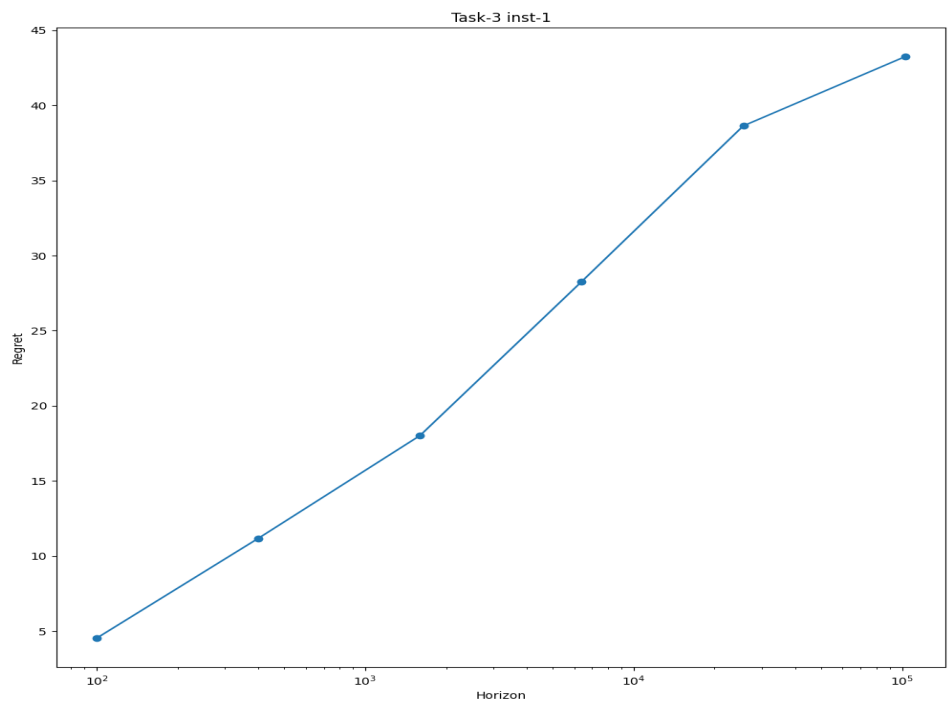
Exploration_bonus is defined as the square root of $c \times \ln(\text{total_pulls}) / \text{pulls_arm}$.

Exploration_bonus increases with c and thus ucb for each arm will increase accordingly.

If scale is too low then empirical mean factor will dominate and if scale is too high then Exploration_bonus factor dominates.

Thus, choosing a proper scale value will help us to choose the optimum arm.

TASK 3



I have used Thompson sampling to implement alg-t3.

For both instances the average regret is less than 50 for a horizon upto 10^5 .

For every arm a , I'm drawing a sample $x_t(a) \sim \text{Beta}(s_t(a) + 1, f_t(a) + 1)$ where $s_t(a)$ are successes and $f_t(a)$ are failures and choosing the arm with maximal $x_t(a)$.

In this algorithm I've pulled each arm for 1 time in the start and then selected an arm depending on the Beta distribution values.

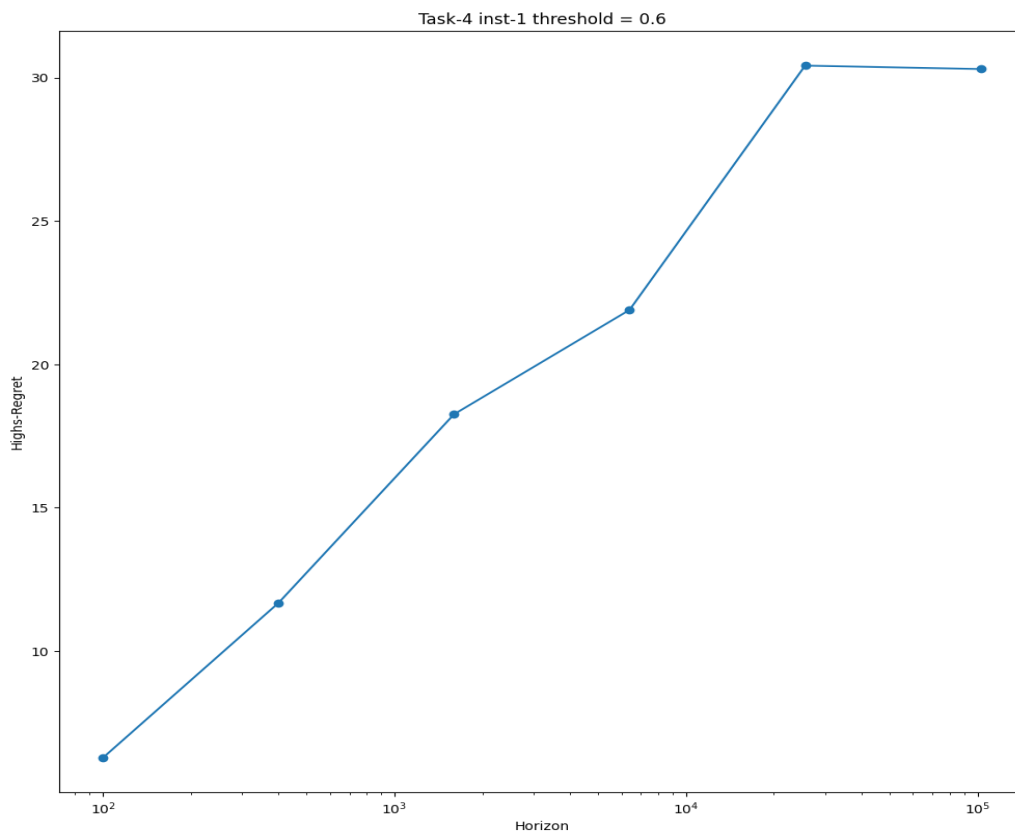
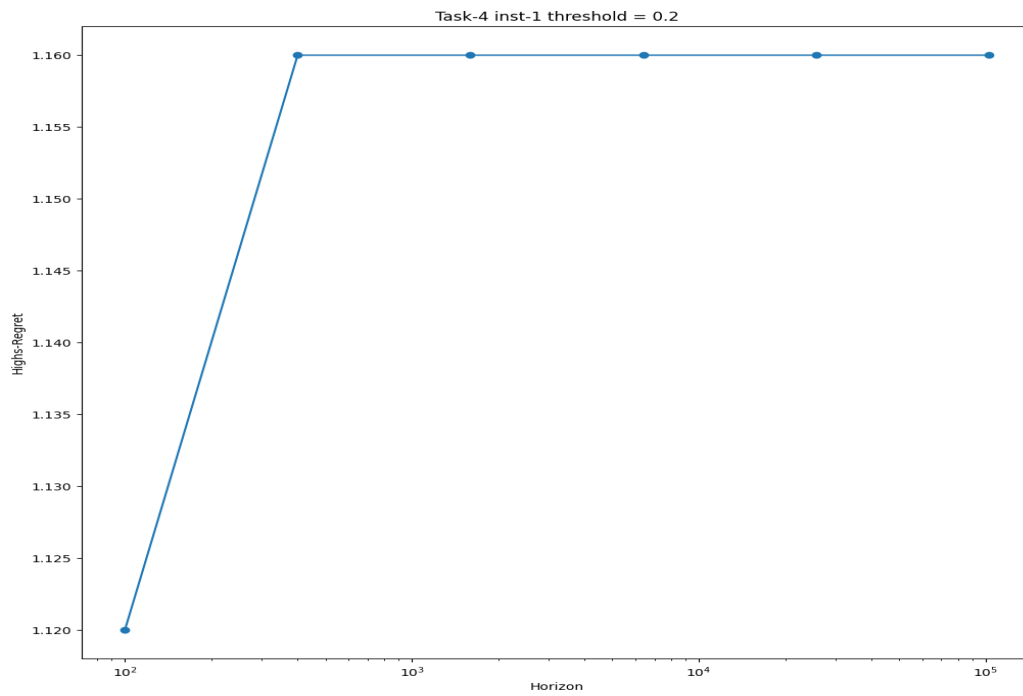
Using `np.random.choice()` I'm generating the reward for each pull.

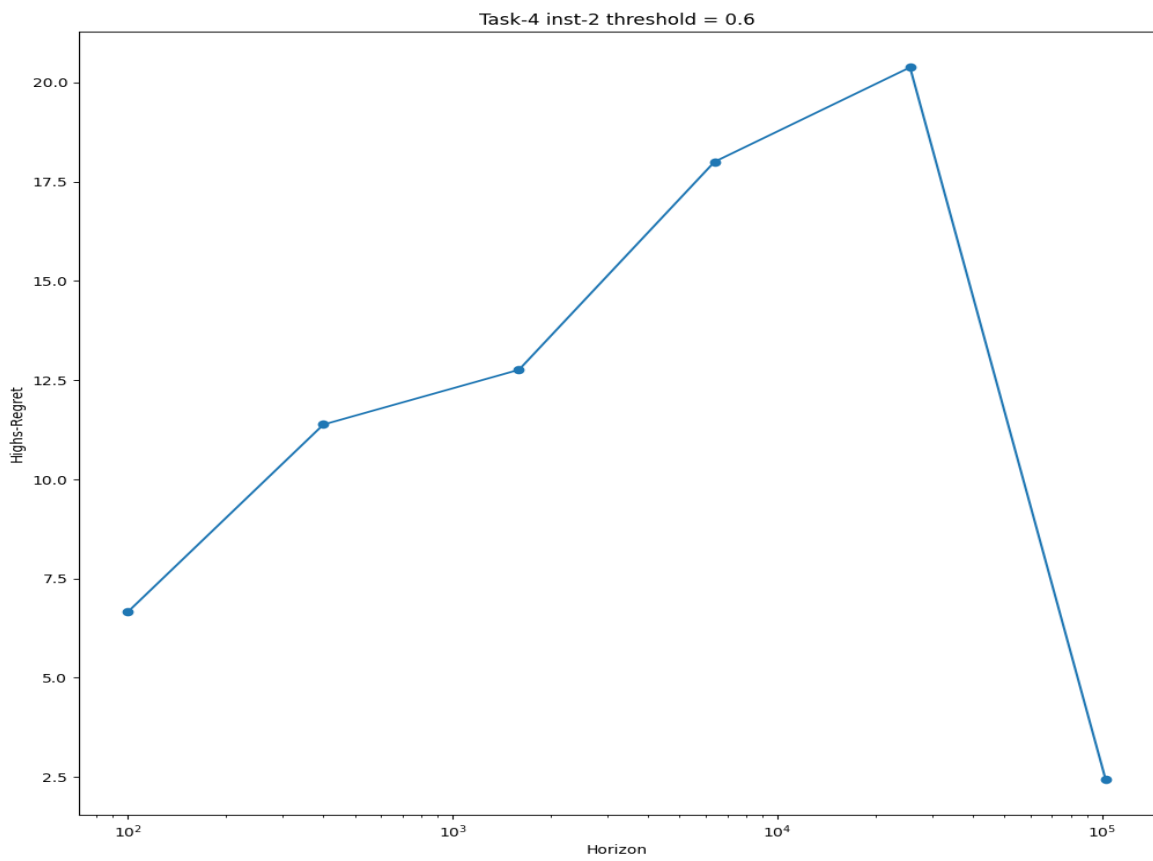
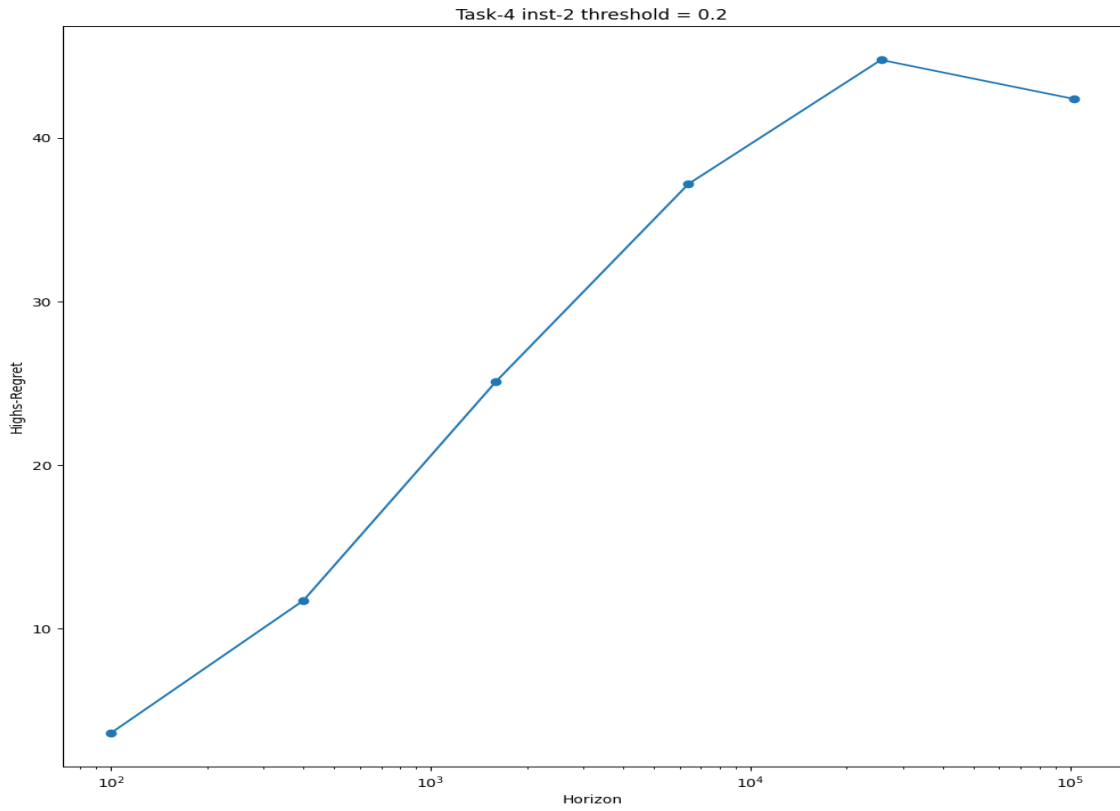
I've calculated regret as follows

1. Getting the maximum of expected rewards obtained from each arm
2. Multiplying it with horizon to get $\text{max_reward} = \text{horizon} \times \max(E(\text{reward}))$
3. Calculating REG as $\text{REG} = \text{max_reward} - \text{REW}$ (cumulative reward)

With the help of Thompson sampling algorithm I'm getting very low regret.

TASK 4





I have used Thompson sampling here as well to implement alg-t4.

In this algorithm I've pulled each arm for 1 time in the start and then selected an arm depending on the Beta distribution values.

Using `np.random.choice()` I'm generating the reward for each pull.

If the reward I get is greater than the threshold then it is considered as '1' else '0'.

I've calculated regret as follows

1. $E(\text{reward})(\text{arm_x}) = p_m + \dots + p_n$ where p_m upto p_n are the probabilities of getting a reward greater than threshold .
2. Getting the maximum of expected rewards obtained from each arm
3. Multiplying it with horizon to get $\text{max_reward} = \text{horizon} \times \max(E(\text{reward}))$
4. Calculating regret as $\text{REG} = \text{max_reward} - \text{REW}$ (cumulative reward)

We can see from the plots that for very large horizons we are getting very less regret.