

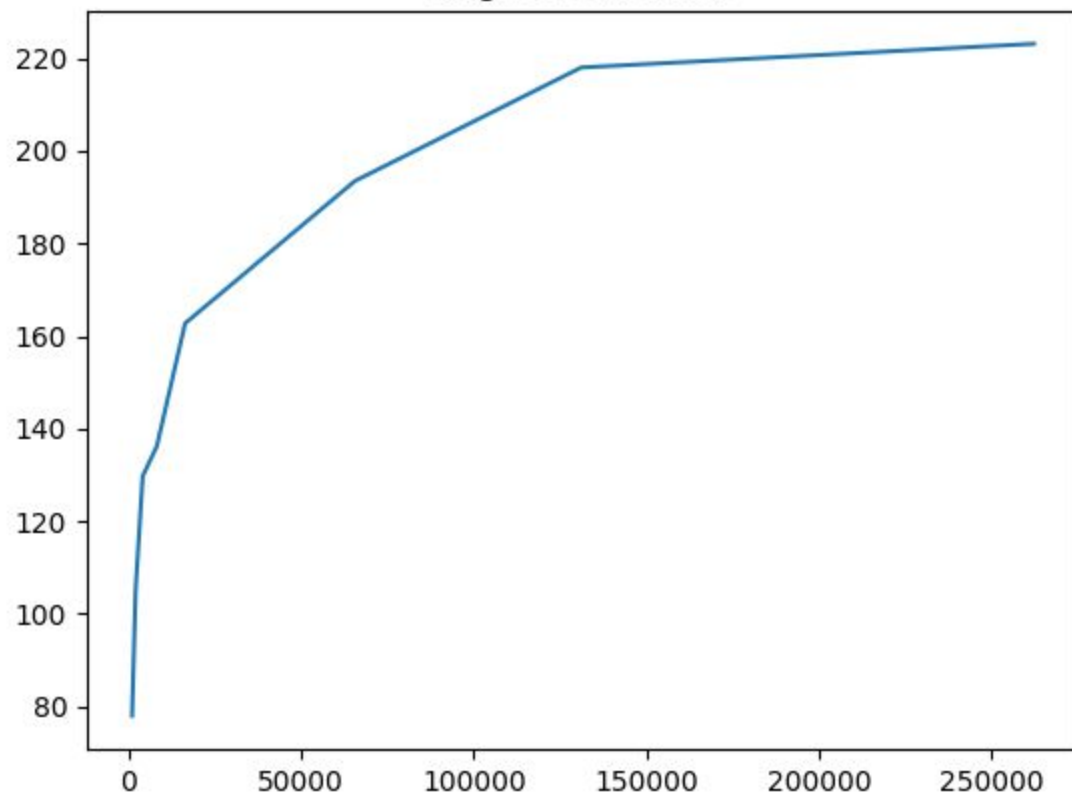
CS747: Assignment_1

Rupansh Parth Kaushik: 200260043

Task 1: KL_UCB

- Def KL(p,q): Calculates $p \cdot \log(p/q) + (1-p) \cdot \log((1-p)/(1-q))$
- Def get_kl_ucb(): Uses bisection method to calculate the solution to the equation $p \cdot \log(p/q) + (1-p) \cdot \log((1-p)/(1-q)) = (\log(t) + 3 \cdot \log(\log(t))) / t$
- Class KL_UCB: Firstly we define the initial variables self.values, self.counts, self.ucb, self.n and self.arms which stores the empirical mean for each arm, the number of times each arm has been pulled, the value of ucb_kl calculated according to the formula given in class slides and which iteration is going on respectively.
- give_pull() : gives the index of the arm which has max value in self.ucb
- get_reward(): updates the values of all the variables.
- The next slide contains the plot from the simulator.py file

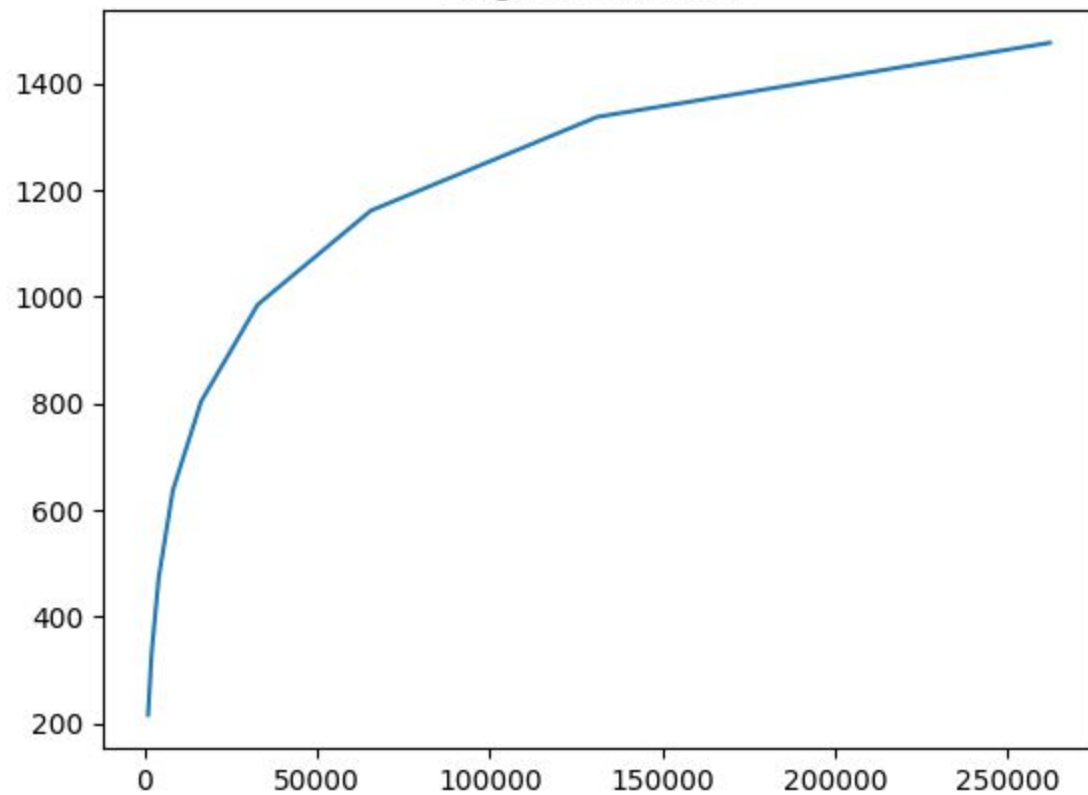
Regret vs Horizon



Task 1: UCB

- `class UCB()`: Firstly we define the initial variables `self.values`, `self.counts`, `self.ucb`, `self.n` and `self.arms` which stores the empirical mean for each arm, the number of times each arm has been pulled, the value of ucb calculated according to the formula given in class slides and which iteration is going on respectively.
- `give_pull()` : gives the index of the arm which has max value in `self.ucb`
- `get_reward()`: updates the values of all the variables.
- The next slide contains the plot from the `simulator.py` file

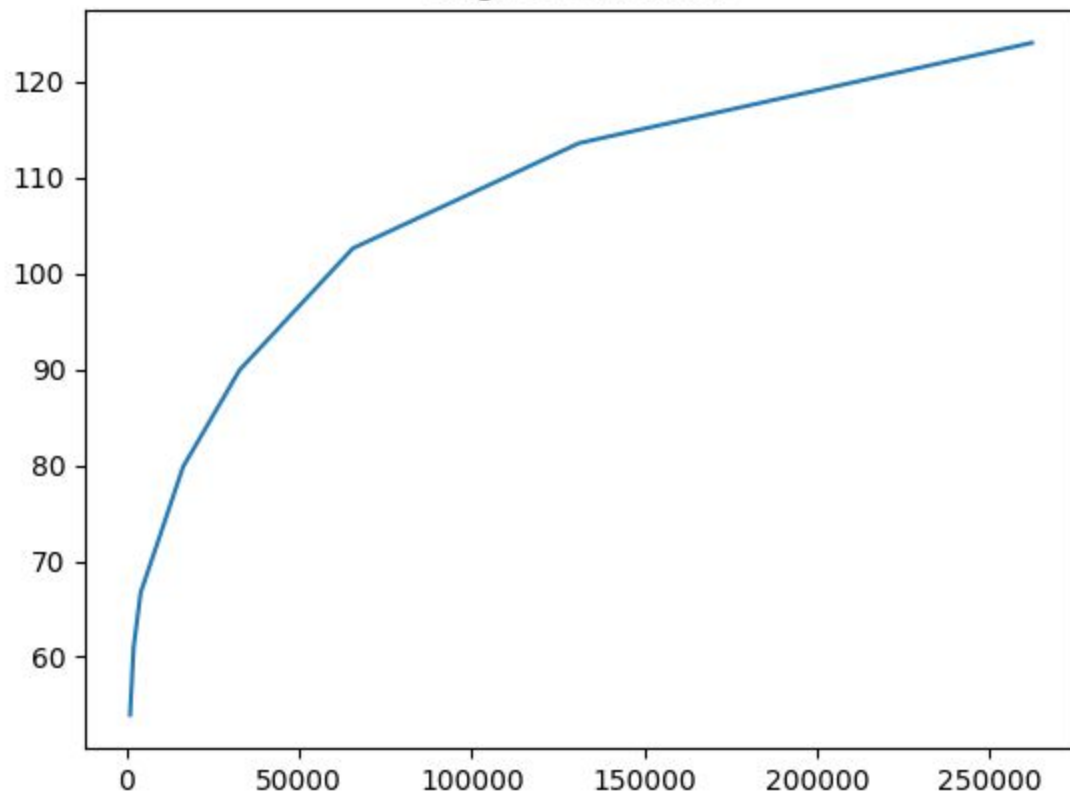
Regret vs Horizon



Task 2: Thompson sampling

- `class Thompson_Sampling()`: Firstly we define the initial variables `self.success`, `self.counts`, `self.beta`, `self.n` and `self.failure` which stores the number of success for each arm, the number of times each arm has been pulled, the value of beta calculated using a numpy beta distribution function, which iteration is going on and number of failures of each arm respectively.
- `give_pull()` : gives the index of the arm which has max value in `self.beta`
- `get_reward()`: updates the values of all the variables.
- The next slide contains the plot from the simulator.py file

Regret vs Horizon



Task 2

My idea for this task was to explore the arms for the first $\text{int}(\text{eps} * \text{Horizon} / \text{batch-size})$ times for some eps and keep updating the beta values for the exploration and then pull a specific number of arms for different batch sizes.

- The pulls will be based on the values of `self.beta` which is same as for `thompson_sampling`.
- And the frequencies depend upon the beta values of each arm. Actually it has been scaled throughout the batch size using the powers of corresponding beta values as reference.
- I was able to get a regret of around 350 for each test case

Task 3:

- First part was exploring the outcomes of the arms for $\frac{1}{4}$ th of the total iterations (Horizon value)
- Second part consist of selecting arms based on the beta values of different arms.