

# CS747: Assignment 1

Aryan Bhosale

210040024

Indian Institute of Technology, Bombay

## I. TASK 1: IMPLEMENTATION OF SAMPLING ALGORITHMS

### **Problem Statement**

*In this first task, you will implement the sampling algorithms: (1) UCB, (2) KL-UCB, and (3) Thompson Sampling. This task is straightforward based on the class lectures. The instructions below tell you about the code you are expected to write.*

In this section, the sampling algorithms Upper Confidence Bound (UCB), Kullback-Leibler Upper Confidence Bound (KL-UCB) and Thompson Sampling have been implemented in code and analysis of regret over different horizons is carried out.

### A. Upper Confidence Bound (UCB) Algorithm

Formula for UCB is:

$$\text{ucb}_a^t = \hat{p}_a^t + \sqrt{\frac{2 \ln(t)}{u_a^t}}$$

where  $\hat{p}_i$  is the empirical mean of rewards from arm  $a$  and  $u_a^t$  is the number of times  $a$  has been sampled at time  $t$ .

The algorithm balances exploration and exploitation to maximize the total reward. The UCB algorithm works by selecting the arm with the highest upper confidence bound. In this implementation, it is ensured that each arm is pulled at least once at random in the beginning after which we can easily update the upper confidence bound of each arm based on the observed rewards. The algorithm then selects the arm with the highest upper confidence bound for exploitation.

1) *Plot:* The exploration bonus ensures that all arms are pulled at some time eventually. Also, the more some arm is pulled, the lesser the exploration bonus for the same. This helps in achieving sub-linear regret. It follows from the plot that regret is  $O(\log(T))$ , as expected from the derivation and analysis of UCB, which proved the regret bound by expanding the summations.

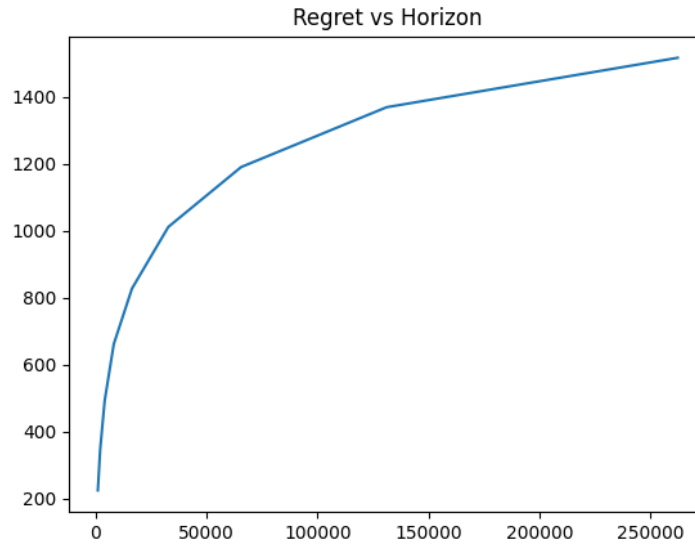


Fig. 1. UCB: Regret v. Horizon

2) *Code Description:* Please refer to coloured annotations for detailed description.

## Variables

Listing 1. Variables and their meanings

```

1 self.num_arms = num_arms # no of arms in the bandits
2 self.curr_timestep = 0 # the current time 't' as defined in the formula
3 self.counts = np.zeros(num_arms) # array to count the no. of times an arm has been pulled
4 self.initial_pulls = np.arange(num_arms) # a shuffled array to ensure we pull all arms atleast
   once in the beginning
5 np.random.shuffle(self.initial_pulls)
6 self.values = np.full((num_arms,), 1e-5) # the empirical mean of rewards
7 self.ucb_arms = np.zeros(num_arms) # array to store upper confidence bound for arms at every
   iteration

```

## Pull Function

Listing 2. get pull function

```

1 ## to ensure all arms are pulled
2 if self.curr_timestep < self.num_arms:
3     return self.initial_pulls[self.curr_timestep]
4 ## choosing the arm with the highest upper confidence bound
5 for i in range(self.num_arms):
6     #formula for upper confidence of an arm at time t
7     self.ucb_arms[i] = self.values[i] + math.sqrt((2*math.log(self.curr_timestep))/self.counts[
        i])
8
9 return np.argmax(self.ucb_arms)

```

## Reward Updation

Listing 3. get reward function

```

1 n = self.counts[arm_index]
2 if n==0:
3     self.values[arm_index] = (self.values[arm_index]+reward)/2
4 else:
5     # taking a running average to store empirical mean of rewards
6     self.values[arm_index] = ((n- 1) / n) * self.values[arm_index] + (1 / n) * reward
7
8 self.curr_timestep += 1
9 self.counts[arm_index] += 1

```

## B. Kullback-Leibler Upper Confidence Bound (KL-UCB) Algorithm

Formula for KL-UCB is:

$$KL-UCB_a^t = \max \left\{ q \in [\hat{p}_a^t, 1] : u_a^t KL(P_i, Q_i(t)) \leq \ln(t) + c \ln(\ln(t)) \right\} \text{ where, } c \geq 3$$

where  $\hat{p}_i$  is the empirical mean of rewards from arm  $a$  and  $u_a^t$  is the number of times  $a$  has been sampled at time  $t$ .

It is an extension of the UCB algorithm that uses the Kullback-Leibler (KL) divergence. The algorithm selects the arm with the highest upper confidence bound, which in this case is calculated using the KL divergence and the number of times the arm has been pulled. The KL-UCB algorithm has been shown to have a logarithmic regret bound, i.e. that the expected loss of the algorithm compared to the best arm decreases logarithmically with the number of rounds. This algorithm tends to be computationally heavy because of the load of calculating the KL-Divergence.

1) *Plot:* It follows from the plot that regret is  $O(\log(T))$ . We know that KL-UCB provides tighter bounds which is normally due to the sub-sub-linear term. The value for regret is significantly lower when compared to UCB. It must be noted that  $c=0$  has been used since that has been proved to be better by the authors of *The KL-UCB Algorithm for Bounded Stochastic Bandits and Beyond*.

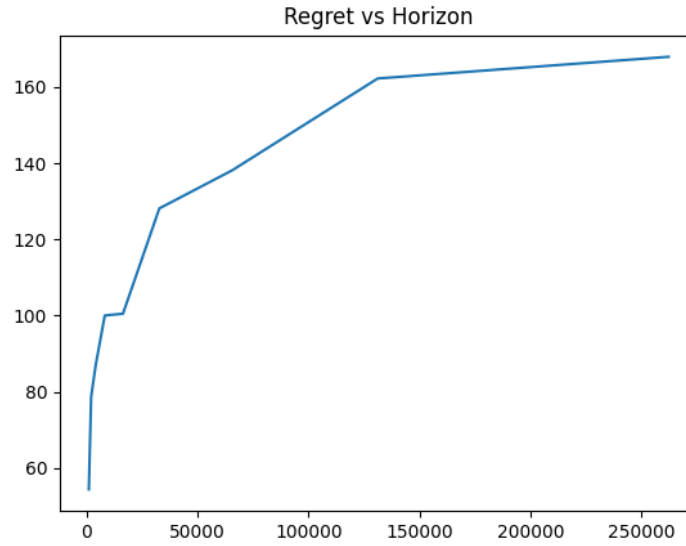


Fig. 2. KL-UCB: Regret v. Horizon

2) *Code Description:* Please refer to coloured annotations for detailed description.

### Variables

Listing 4. Variables and their meanings (same as UCB)

```

1 self.num_arms = num_arms # no of arms in the bandits
2 self.curr_timestep = 0 # the current time 't' as defined in the formula
3 self.counts = np.zeros(num_arms) # array to count the no. of times an arm has been pulled
4 self.initial_pulls = np.arange(num_arms) # a shuffled array to ensure we pull all arms atleast
   once in the beginning
5 np.random.shuffle(self.initial_pulls)
6 self.values = np.full((num_arms, ), 1e-5) # the empirical mean of rewards
7 self.ucb_arms = np.zeros(num_arms) # array to store upper confidence bound for arms at every
   iteration

```

### Pull Function

Listing 5. get pull function

```

1  ## to ensure all arms are pulled
2  if self.curr_timestep < self.num_arms:
3      return self.initial_pulls[self.curr_timestep]
4  ## choosing the arm with the highest upper confidence bound calculated using KL divergence
5  ## Binary search has been used below to calculate KL-UCB
6  for i in range(self.num_arms):
7      p_a_t = self.values[i]
8      lower = p_a_t
9      upper = 1
10     tolerance = 0.01
11
12     while (upper - lower)>tolerance:
13         mid = (lower+upper)*0.5
14         kl_term = kl_div(p_a_t, mid)
15         if kl_term < math.log(self.curr_timestep)/self.counts[i] :
16             lower = mid
17         else:
18             upper = mid
19     self.kl_ucb_arms[i] = (lower+upper)*0.5
20
21 return np.argmax(self.kl_ucb_arms)

```

## Reward Updation

Listing 6. get reward function (same as UCB)

```

1  n = self.counts[arm_index]
2  if n==0:
3      self.values[arm_index] = (self.values[arm_index]+reward)/2
4  else:
5      # taking a running average to store empirical mean of rewards
6      self.values[arm_index] = ((n- 1) / n) * self.values[arm_index] + (1 / n) * reward
7
8  self.curr_timestep += 1
9  self.counts[arm_index] += 1

```

## KL Divergence Function

Listing 7. kldiv

```

1  ##function to calculate the KL divergence between a given x and y
2  def kl_div(x, y):
3      if x == 0:
4          return math.log(1/(1-y))
5      elif x==1:
6          return math.log(1/y)
7      term1 = x*(math.log(x/y))
8      term2 = (1-x)*math.log((1-x)/(1-y))
9      return term1 + term2

```

## C. Thompson Sampling Algorithm

According to this algorithm, for every arm we draw a sample from the corresponding beta distribution  $x_a^t$  which is defined as follows.

If at time  $t$ , an arm has succeeded  $s_a^t$  times and failed  $f_a^t$  times the beta distribution is defined as:  $\beta(s_a^t + 1, f_a^t + 1)$  where  $\alpha$  is  $s_a^t + 1$  and  $\beta$  is  $f_a^t + 1$  which implies  $x_a^t$  for every arm is:  $x_a^t \sim \text{Beta}(s_a^t + 1, f_a^t + 1)$ . We pull the arm that has the highest  $x_a^t$ .

1) *Plot:* We have achieved sub-linear regret. Additionally this algorithm is very easy to implement. It is known to work well in practice. The random sampling using Beta distribution accounts for the empirical mean and associated uncertainty with it. We add one so that the coefficients are always positive irrespective of the value of successes and failures. Thompson sampling has clearly minimised regret better than both the UCB algorithms as regret is curtailed to below 120. We can also see that regret is  $O(\log(T))$ .



Fig. 3. Thompson: Regret v. Horizon

2) *Code Description:* Please refer to coloured annotations for detailed description.  
Variables

Listing 8. Variables and their meanings

```
1 self.num_arms = num_arms # no of arms in the bandits
2 self.successes = np.zeros(num_arms) #no. of times a pull for an arm was successful
3 self.fails = np.zeros(num_arms) #no. of times a pull for an arm was a failure
4 self.sampler_vals = np.zeros(num_arms) # array to store the sampled value for arms at every
   iteration
```

### Pull Function

Listing 9. get pull function

```
1 ## choosing the arm with the highest sampled value from its beta distribution
2 for i in range(self.num_arms):
3     alpha = self.successes[i] + 1
4     beta = self.fails[i] + 1
5     self.sampler_vals[i] = np.random.beta(alpha, beta)
6 return np.argmax(self.sampler_vals)
```

### Reward Updation

Listing 10. get reward function

```
1 #updating the no. of failures and successes
2 if(reward > 0):
3     self.successes[arm_index] += 1
4 else:
5     self.fails[arm_index] += 1
```

## II. TASK 2 A

### Problem Statement

This task explores the effect of the difference between the means of arms on the regret accumulated by the UCB algorithm. For this task, take two-armed bandit instances (with means  $[p_1, p_2]$ ) with the higher mean arm fixed (say,  $p_1$ ), and vary the other arm's mean ( $p_2$ ) from 0 to  $p_1$ . For the assignment,  $p_1 = 0.9$ , and  $p_2$  varies from 0 to 0.9 (both inclusive) in steps of 0.05. Do this for a horizon of 30000. Plot

the variation of regret with  $p_2$ . Clearly state your observations from the plot, and explain the reason(s) behind the observations.

1) *Plot:* When  $p_2$  is closer to  $p_1$  (i.e. around 0.9), the regret is higher since it's more challenging to distinguish the arms. When  $p_2$  is significantly lower than  $p_1$ , the regret is low because the UCB algorithm identifies the better arm more quickly and exploits it. The plot illustrates the balance between exploration (choosing arms with uncertain outcomes) and exploitation (choosing the arm with the highest estimated mean) in the UCB algorithm. The algorithm initially explores more when  $p_2$  is close to  $p_1$  but shifts towards exploitation as the difference between  $p_1$  and  $p_2$  becomes more significant.

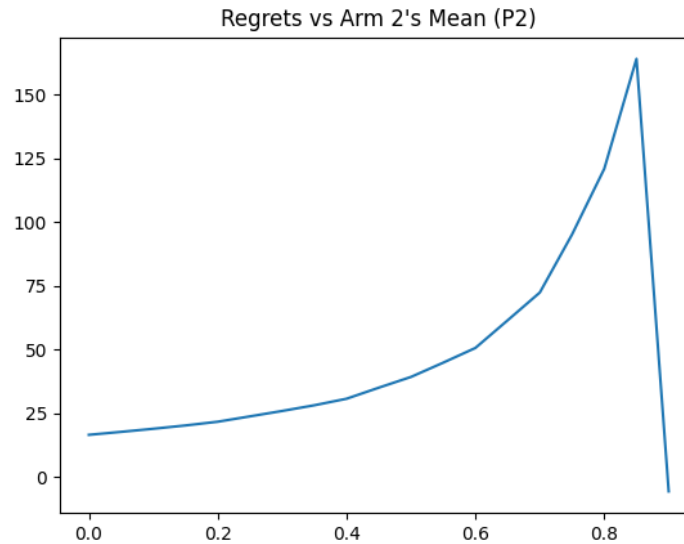


Fig. 4. Task 2A: Regret v. Mean of Arm 2 ( $p_2$ )

Listing 11. get pull function

```

1 #creating the array of arm means as described
2 task2Ap2s = np.arange(0, 0.901, 0.05)
3 task2Ap1s = np.full((task2Ap2s.shape), 0.9)
4
5 #calling the function to execute UCB algorithm as written in the code in Task 1
6 regrets_ucb_2A = task2(UCB, 30000, task2Ap1s, task2Ap2s)

```

### III. TASK 2 B

#### **Problem Statement**

This task involves comparing the behaviour of UCB and KL-UCB algorithms on two-armed bandits. The task is to compare the effect of the value of the means on the algorithms while keeping the difference between the means fixed. For this task, take  $\Delta = p_1 - p_2 = 0.1$ . Again, vary  $p_2$  from 0 to 0.9 (both inclusive) in steps of 0.05 for a horizon of 30000. Plot the variation of regret for both algorithms on different plots with respect to the instances (you may plot with respect to either  $p_1$  or  $p_2$ , but clearly label that in your plots). State and compare your observations for the variations of regret for each algorithm and explain the reason behind your observations.

## Part 1: UCB

1) *Plot:* The plot for UCB algorithm has been shown in the figure.

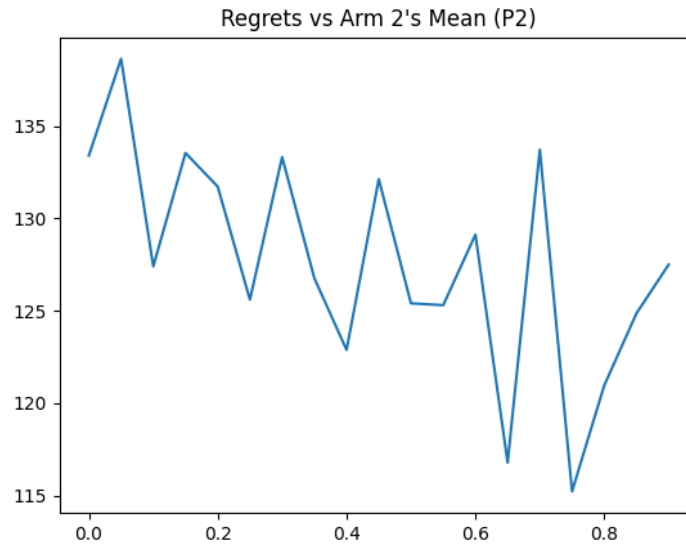


Fig. 5. Task 2B UCB: Regret v. Mean of Arm 2 (p2)

## 2) Code Description: Algorithm

Listing 12. algorithm calling function

```
1 #creating the array of arm means as described
2 task2Bp2s = task2Ap2s
3 task2Bp1s = np.add(task2Bp2s, 0.1)
4
5 #calling the function to execute UCB algorithm as written in the code in Task 1
6 regrets_ucb_2B = task2(UCB, 30000, task2Bp1s, task2Bp2s)
```

## Plotter Function

Listing 13. plotter function

```
1 #x, y are the respective axes which get p2 and the regrets respectively
2 def plotter(x, y, what_this_is):
3     plt.plot(x, y)
4     title = "Regrets vs Arm 2's Mean (P2)"
5     plt.title(title)
6     #the plot is then saved as a .png file
7     plt.savefig("task1-{}-{}.png".format(what_this_is, time.strftime("%Y%m%d-%H%M%S")))
8     plt.clf()
```

## Part 2: KL-UCB

### A. Plot

The plot for KL-UCB algorithm has been shown in the figure.

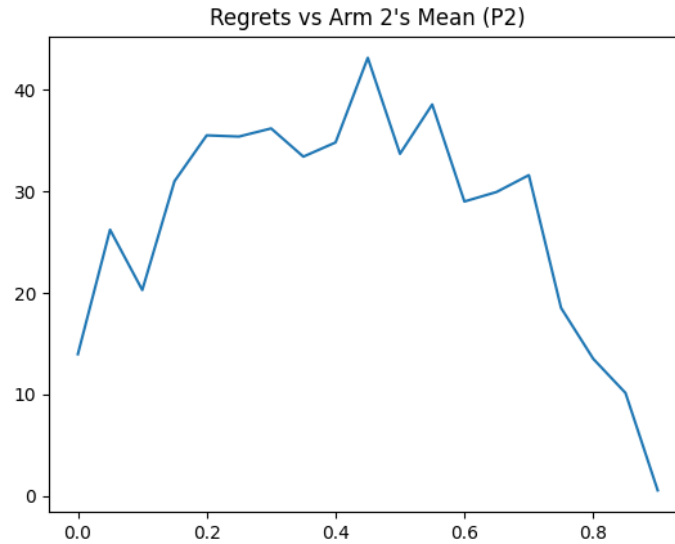


Fig. 6. Task 2B KL-UCB: Regret v. Mean of Arm 2 (p2)

### 1) Code Description: Algorithm

Listing 14. algorithm calling function

```

1 #creating the array of arm means as described
2 task2Bp2s = task2Ap2s
3 task2Bp1s = np.add(task2Bp2s, 0.1)
4
5 #calling the function to execute UCB algorithm as written in the code in Task 1
6 #horizon = 30,000
7 regrets_kl_ucb_2B = task2(KL_UCB, 30000, task2Bp1s, task2Bp2s)

```

### Interpretation

Graph 1 - UCB Graph 2 - KL-UCB As expected KL-UCB is better as minimizing regret since it provides tighter bounds. Additionally when both arms have high values of means, we see that it does not matter which arm is chosen which results in a lesser value of regret. The peak at  $p2 = 0.5$  (approximately) can be explained as follows. In this scenario, the algorithm may find it challenging to distinguish between them, leading to more exploration and higher regret.



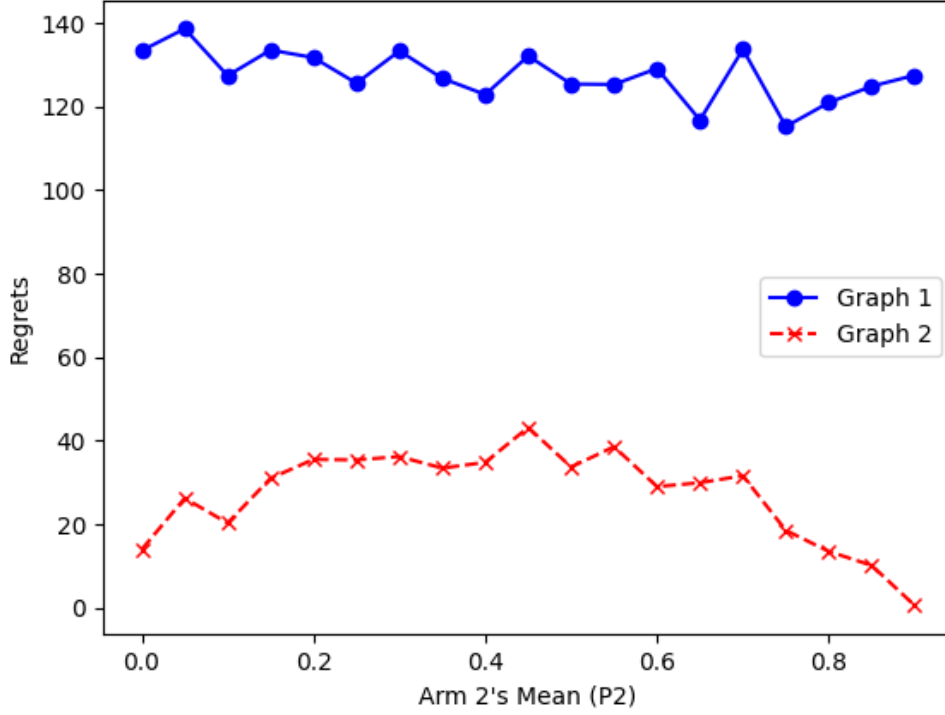


Fig. 7. Task 2B Comparative: Regret v. Mean of Arm 2 (p2)

#### IV. TASK 3

##### **Problem Statement**

*This task involves dealing with a bandit instance where your pulls are no longer guaranteed to be successful and have a probability of giving faulty outputs. When you pull an arm, it has a certain known probability of giving the correct output of the arm, otherwise it returns a 0 or 1 uniformly at random. Your task is to come up with a good algorithm to maximise the reward for this faulty bandit setting.*

##### **All Tried Solutions**

The algorithm used is a modification of the UCB algorithm. The only change is in the reward updation function. Since we know the probability with which we get a faulty output, at each iteration, we will decide whether we are considering the obtained reward as the output of the bandit or a random result. If it is an actual result, we update the arm means and counts. Otherwise we just increase the no. of timesteps so that we can have higher exploration bonus. This algorithm was also tried by using KL-UCB as the bound based on which the decision is taken however that lead to higher computational load and did not improve rewards by much which is why that method was discarded.

##### **1) Code Description: Algorithm**

Listing 15. algorith calling function

```

1 #deciding whether this reward is due to a fault or an actual pull of an arm
2 fault_or_not = np.random.choice([0,1], p=[self.fault, self.not_fault])
3 #0==fault, 1==not fault
4 if fault_or_not:
5     n = self.counts[arm_index]
6     if n==0:

```

```

7         self.values[arm_index] = (self.values[arm_index]+reward)/2
8     else:
9         self.values[arm_index] = ((n- 1) / n) * self.values[arm_index] + (1 / n) * reward
10        self.good_ones += 1
11        self.counts[arm_index] += 1
12    #updating the no. of timesteps irrespective of which case were in
13    self.curr_timestep += 1

```

### *Best Solution from All Tried Solutions*

When Thompson Sampling was run without any modifications, we find that it outperforms all the above solutions. It seems to be immune to this faulty bandit problem. This is why Thompson Sampling has been used as the solution for this Task instead of the above listed methods. Thompson Sampling has been described in Task 1 Section C.

## V. TASK 4: MULTI-MULTI-ARM BANDIT

### *Problem Statement*

*This task involves dealing with two bandit instances at once! In this task, whenever you specify an arm index to pull, one of two given bandit instances is chosen uniformly at random, and the arm corresponding to your provided index is pulled (both instances have equal number of arms). Once an arm is pulled, the environment returns the reward obtained along with which bandit instance was chosen for that pull. Your task is to come up with a good algorithm to maximise the reward for this multi-multi-armed bandit setting.*

### *All Tried Solutions*

In this we use a different modification of UCB algorithm. Since we do not know the set from which the arm at that arm, we maintain UCB values for all arms of both sets, take an average of ucb values for each arm index and then pick the arm for which the averaged UCB value is maximum. This was tried with KL-UCB but again, it took longer to run and did not improve performance by much which is why it was discarded. A variation of Thompson Sampling was also tried where I averaged the sampled values from the beta distribution for the 2 sets but that gave way worse results because averaging in that case is clearly useless.

#### *1) Code Description: Variables*

Listing 16. Variables

```

1  ## Twice the number of variables have been defined as in Vanilla UCB since we are keeping track
   of both sets separately
2  self.curr_timestep = 0
3  #to check if all arms have been pulled atleast once
4  self.once_0 = 0
5  self.once_1 = 0
6  #no. of pulls for an arm
7  self.counts_0 = np.full((num_arms,), 1e-5)
8  self.counts_1 = np.full((num_arms,), 1e-5)
9  self.initial_pulls = np.arange(num_arms)
10 np.random.shuffle(self.initial_pulls)
11 #empirical mean of rewards
12 self.values_0 = np.full((num_arms,), 1e-5)
13 self.values_1 = np.full((num_arms,), 1e-5)
14 self.ucb_arms_0 = np.zeros(num_arms)
15 self.ucb_arms_1 = np.zeros(num_arms)
16 self.actual_ucbs = np.zeros(num_arms)

```

#### **Pull Function**

Listing 17. Variables

```

1  ##ensuring that all arms are pulled at least once
2  if self.once_0 + self.once_1 < 2*self.num_arms:
3      if self.once_0 < self.once_1:
4          return self.initial_pulls[self.once_0]
5      else:
6          return self.initial_pulls[self.once_1]
7  ##taking average of UCB values of both arms of same index
8  for i in range(self.num_arms):
9      self.ucb_arms_0[i] = self.values_0[i] + math.sqrt((2*math.log(self.curr_timestep))/self.
10         counts_0[i])
11      self.ucb_arms_1[i] = self.values_1[i] + math.sqrt((2*math.log(self.curr_timestep))/self.
12         counts_1[i])
13      self.actual_ucbs[i] = (self.ucb_arms_0[i] + self.ucb_arms_1[i])*0.5
14  return np.argmax(self.actual_ucbs)

```

## Reward Updation

Listing 18. Variables

```

1  ##if arm from the first set is pulled
2  if set_pulled == 0:
3      n = self.counts_0[arm_index]
4      # if n==1e-5:
5      #     self.values_0[arm_index] = (self.values_0[arm_index]+reward)/2
6      # else:
7      #     self.values_0[arm_index] = ((n- 1) / n) * self.values_0[arm_index] + (1 / n) *
8          reward
9      self.values_0[arm_index] = ((n- 1) / n) * self.values_0[arm_index] + (1 / n) * reward
10     self.once_0 += 1
11     self.counts_0[arm_index] += 1
12 ##if arm from the second set is pulled
13 else:
14     n = self.counts_1[arm_index]
15     # if n==1e-5:
16     #     self.values_1[arm_index] = (self.values_1[arm_index]+reward)/2
17     # else:
18     #     self.values_1[arm_index] = ((n- 1) / n) * self.values_1[arm_index] + (1 / n) *
19         reward
20     self.values_1[arm_index] = ((n- 1) / n) * self.values_1[arm_index] + (1 / n) * reward
21     self.once_1 += 1
22     self.counts_1[arm_index] += 1
23 self.curr_timestep += 1

```

## Best Solution from All Tried Solutions

When Thompson Sampling was run without any modifications, we find that it outperforms all the above solutions. It seems to be immune to this faulty bandit problem. This is why Thompson Sampling has been used as the solution for this Task instead of the above listed methods. Thompson Sampling has been described in Task 1 Section C.