

1 Task 1

1.1 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.

Read task1.py. It contains a sample implementation of epsilon-greedy for you to understand the Algorithm class. You have to edit the `__init__` function to create any state variables your algorithm needs, and implement `give_pull` and `get_reward`. `give_pull` is called whenever it is the algorithm's decision to pick an arm and it must return the index of the arm your algorithm wishes to pull (lying in 0, 1, ... `self.num_arms-1`). `get_reward` is called whenever the environment wants to give feedback to the algorithm: your code can use this feedback to update the data structures maintained by your agent. It will be provided the `arm_index` of the arm and the reward seen (0/1). Note that the `arm_index` will be the same as the one returned by the `give_pull` function called earlier. For more clarity, refer to `single_sim` function in `simulator.py`.

Once done with the implementations, you can run `simulator.py` to see the regrets over different horizons. Save the generated plot and add it to your report, with suitable commentary (ideally 4-5 lines describing what you see, what issues you faced, any surprising patterns). You may also run `autograder.py` to evaluate your algorithms on the provided test instances.

1.2 UCB

In UCB we need to pull the arm which has the highest empirical mean. But before doing that we should pull every arm once and calculate the empirical mean of it. We can do it by storing number of pulls we have made till now in a variable. We will be storing empirical mean array, array which stores number of times an arm got pulled and an array to store the rewards it got us.

```
class UCB(Algorithm):
    def __init__(self, num_arms, horizon):
        super().__init__(num_arms, horizon)
        # You can add any other variables you need here
        # START EDITING HERE
        self.counts=np.zeros(num_arms) #to store the number of times each arm is pulled
        self.ucbs=np.zeros(num_arms) #to store ucb of each arm
        self.rews=np.zeros(num_arms) #to store total rewards of each arm
        self.means=np.zeros(num_arms) #to store empirical means of each arms
        self.pulls=0
        self.zero_ind=0
        # END EDITING HERE
```

This is the snippet of my init function for UCB.

Let us see the `give_pull` function now. At first it pulls every arm once and from then select the arm with maximum empirical mean.

```
114 def give_pull(self):
115     # START EDITING HERE
116     if self.zero_ind<self.num_arms:
117         return self.zero_ind
118     else:
119         return np.argmax(self.ucbs)
120     # END EDITING HERE
```

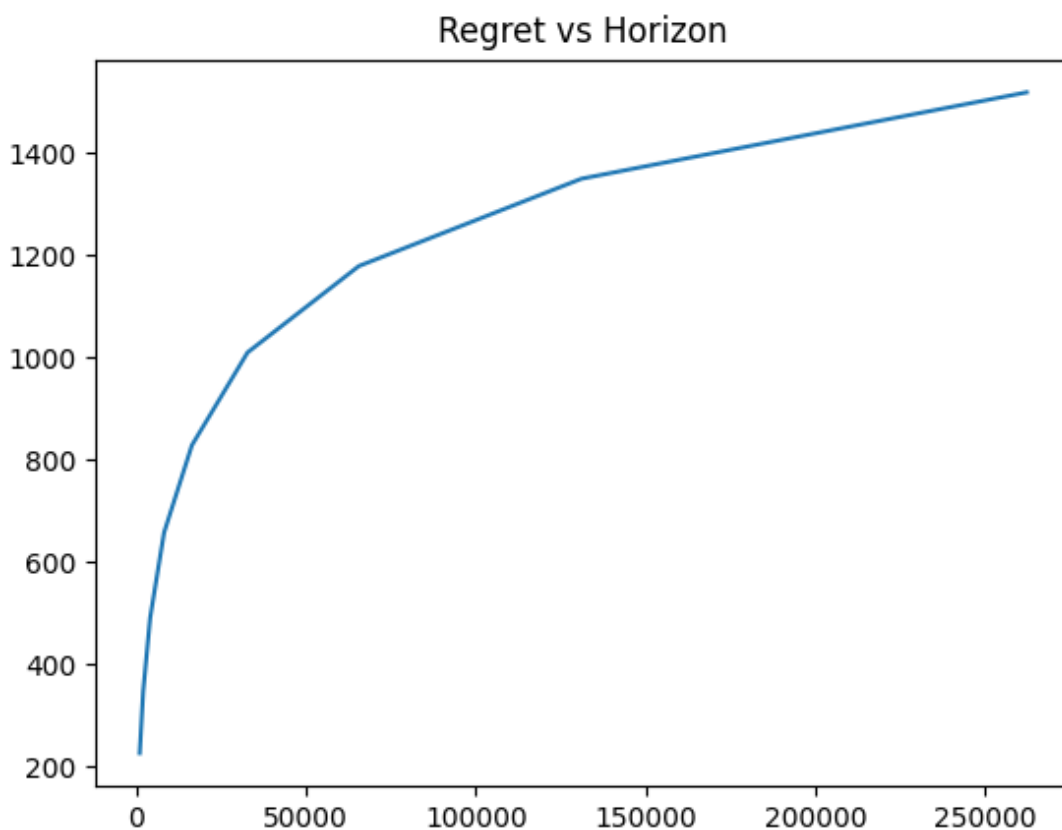
This is the snippet of my `get_pull` function for UCB.

Let us see the `get_reward` function now. Here we need to update the arrays which stores the means and counts. i have divided in to three conditions based on my requirements. Each condition is based on number of pulls we have done.

```
122     def get_reward(self, arm_index, reward):
123         # START EDITING HERE
124         if self.zero_ind < self.num_arms - 1:
125             self.rews[arm_index] += reward
126             self.counts[arm_index] += 1
127             self.pulls += 1
128             self.zero_ind += 1
129         elif self.zero_ind == self.num_arms - 1:
130             self.rews[arm_index] += reward
131             self.counts[arm_index] += 1
132             self.pulls += 1
133             self.zero_ind += 1
134             self.means = np.divide(self.rews, self.counts)
135             x = np.sqrt(2 * math.log(self.pulls) / self.counts)
136             self.ucbs = self.means + x
137         else:
138             self.counts[arm_index] += 1
139             self.rews[arm_index] += reward
140             self.pulls += 1
141             self.means = np.divide(self.rews, self.counts)
142             x = np.sqrt(2 * math.log(self.pulls) / self.counts)
143             self.ucbs = self.means + x
```

This is the snippet of my `get_reward` function for UCB.

Let us see the plot of Regret vs Horizon we got by simulating UCB.



This is the plot of we get for UCB.

As we can see the plot is logarithmic with respect to N as expected

1.3 KL UCB

KL-UCB is similar to UCB but instead of choosing the arm with highest empirical mean we choose the arm with highest q (Upper Confidence Bound). We calculate q by using kl function and using an algorithm similar to binary search. Let us see the init function first. In init I am storing the counts, rewards, empirical means, kl_bounds (value of q) values of every arm.

```
147 class KL_UCB(Algorithm):
148     def __init__(self, num_arms, horizon):
149         super().__init__(num_arms, horizon)
150         # You can add any other variables you need here
151         # START EDITING HERE
152
153         self.counts=np.zeros(num_arms) #to store the number of times each arm is pulled
154         self.klucbs=np.zeros(num_arms) #to store ucb of each arm
155         self.rews=np.zeros(num_arms) #to store total rewards of each arm
156         self.means=np.zeros(num_arms) #to store empirical means of each arms
157         selfpulls=0
158         self.zero_ind=0
159
```

This is the snippet of my init function for KL-UCB.

Give pull function in here is same as UCB except that we chose arm based on KL bounds rather than on empirical means and even the get reward function is similar to that of UCB with some slight modifications.

```
164 def give_pull(self):
165     # START EDITING HERE
166     if self.zero_ind<self.num_arms:
167         #print(self.zero_ind)
168         return self.zero_ind
169     else:
170         #print(np.argmax(self.klucbs))
171         return np.argmax(self.klucbs)
172     # END EDITING HERE
```

This is the snippet of my give_pull function for KL-UCB.

As I have mentioned it, I used a function which is very similar to binary search but not exactly same. Let us take a look at that function too. It returns the q value for a given arm and based on it we can select an arm to pull.

```
76 def q_calc(pa,ua,rhs):
77     if pa==1:
78         return 1
79     i=np.float64(pa)
80     j=np.float64(0.999)
81     rhs=rhs/ua
82     while(j-i>=0.0001):
83         z=(i+j)/2
84         if(kl(pa,z)<rhs):
85             i=z+0.0001
86         else:
87             j=z-0.0001
88     return (i+j)/2
89
90
```

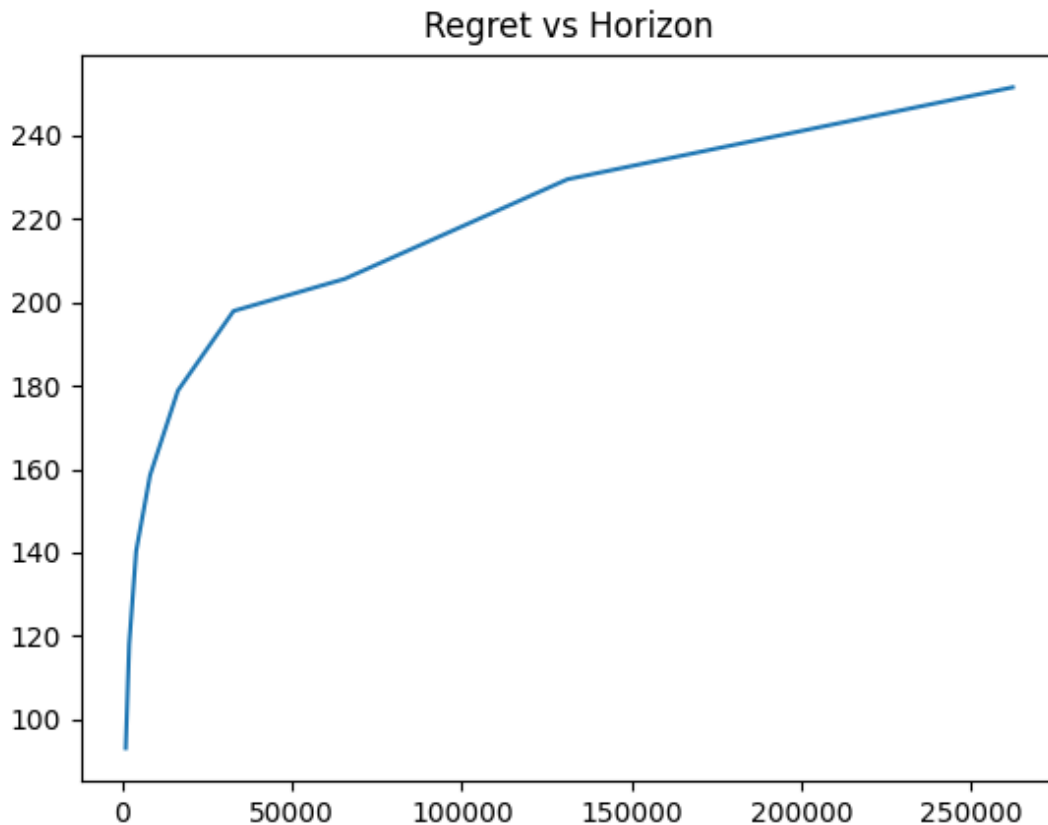
This is the snippet of my binary search function for KL-UCB.

Here is the formula to calculate upper confidence bound of an arm.

$$ucb-kl_a^t = \max\{q \in [\hat{p}_a^t, 1] \text{ such that } u_a^t KL(\hat{p}_a^t, q) \leq \ln(t) + c \ln(\ln(t))\}, \text{ where } c \geq 3.$$

KL-UCB algorithm: at step t , pull $\operatorname{argmax}_a ucb-kl_a^t$.

Let us look at the plot of regret vs horizon we are getting for KL-UCB. It is also logarithmic as expected.



This is plot of Regret vs Horizon function for KL-UCB.

1.4 Thompson Sampling

This algorithm is quite different from UCB, KL-UCB. In this algorithm we store number of successes and failures each arm have. Then we construct a beta distribution on these two values for each arm. Then we take a random value from each distribution and we choose the arm with the highest random value. It is a non-deterministic algorithm while the other two are deterministic. Let us look at the init function in thompson sampling.

```
203
204 class Thompson_Sampling(Algorithm):
205     def __init__(self, num_arms, horizon):
206         super().__init__(num_arms, horizon)
207         # You can add any other variables you need here
208         # START EDITING HERE
209         self.succs=np.ones(num_arms)
210         self.fails=np.ones(num_arms)
211         # END EDITING HERE
```

This is the snippet of my init function for thompson sampling.

The pull function is quite simple in here. We just need to select the arm which gets a higher random value from beta distribution we have. Let us see the code snippet of pull.

```

212
213 def give_pull(self):
214     # START EDITING HERE
215     arr=np.zeros(self.num_arms)
216     for i in range(self.num_arms):
217         arr[i]=np.random.beta(self.succs[i],self.fails[i])
218     return np.argmax(arr)
219     # END EDITING HERE
220

```

This is the snippet of my give_pull function for thompson sampling.

The get_reward function is pretty straightforward. We just need to change the number of successes or failures of the arm we pulled. Let us see the code snippet for it.

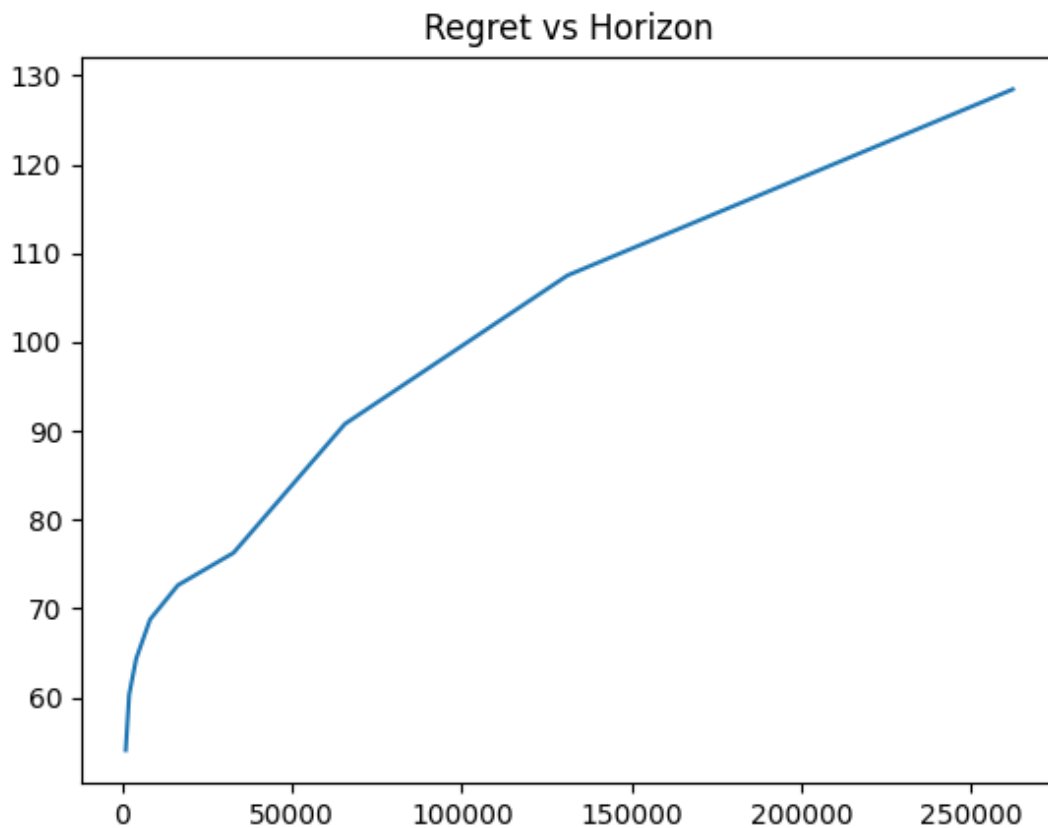
```

220
221 def get_reward(self, arm_index, reward):
222     # START EDITING HERE
223     if reward==1:
224         self.succs[arm_index]+=1
225     else:
226         self.fails[arm_index]+=1
227     # END EDITING HERE
228

```

This is the snippet of my get_reward function for thompson sampling.

Finally let us see the plot of Regret vs Horizon of Thomson's sampling



This is the plot of Regret vs Horizon function for thompson sampling.

2 Task2

2.1 Problem statement

We describe the problem statement that was briefly discussed completely now. The algorithm is given the number of arms of the bernoulli bandit `num_arms`, a horizon and a `batch_size`. The `give_pull` function will be called `horizon/batch_size` times (which can be assumed to be an integer). In every call, it must return the next `batch_size` number of pulls the algorithm wishes to make. The function must do so in a specific format: it has to return two lists, one containing the arm indices that it wishes to pull, and the other containing the number of times each of those indices must be pulled. For example, in a 10-armed bandit instance with `batch_size` 20, a possible return of the `give_pull` function could be `([2, 4, 9], [10, 4, 6])`. Note that your function should generalize to arbitrary `batch_size`s, as long as the given `batch_size` is a factor of the horizon (the `batch_size` could be just 1, or it could also be of the order of the horizon). The autograder/simulator will proceed and pull these arms according to their respective counts, and then provide feedback to the `get_reward` function. The feedback is provided as a dict, where the keys are the arm indices, and the rewards are a numpy array of 0s and 1s that were seen. So a possible input (`arm_rewards`) to the `get_reward` function for the above batch pull could be 2: `np.array([1, 1, 1, 0, 1, 1, 0, 1, 0, 1])`, 4: `np.array([1, 1, 0, 0])`, 9: `np.array([0, 1, 0, 1, 0, 0])`. Again, you can read `single_batch_sim` for more clarity. The regret is calculated over all the pulls over the horizon. Once done with your implementation, you can run `simulator.py` to see the regrets for a fixed horizon over different batch sizes. Save the generated plot and add it to your report, with apt captioning. Take 4-5 lines (or more) to explain the trend you see, and justify your choice of the distribution of pulls in a given `batch_size`. You may also run `autograder.py` to evaluate your algorithms on the provided test instances.

2.2 batch sampling

We can successfully do the batch sampling using the thompson's algorithm. Instead of choosing just one arm we should select batch. So instead of doing one draw from beta distributions we will do multiple draws and store the arm with maximum random values in the array. Syntax is very similar to thompson. Let us look at the code of it.

```
31 class AlgorithmBatched:
32     def __init__(self, num_arms, horizon, batch_size):
33         self.num_arms = num_arms
34         self.horizon = horizon
35         self.batch_size = batch_size
36         assert self.horizon % self.batch_size == 0, "Horizon must be a multiple of batch size"
37         # START EDITING HERE
38         # Add any other variables you need here
39         self.succs=np.ones(num_arms)
40         self.fail= np.ones(num_arms)
41         # END EDITING HERE
42
43     def give_pull(self):
44         # START EDITING HERE
45         l1=[]
46         l2=[]
47         init=[]
48         arr=np.zeros(self.num_arms)
49         for k in range(self.batch_size):
50             for i in range(self.num_arms):
51                 arr[i]=np.random.beta(self.succs[i],self.fail[i])
52             init.append(np.argmax(arr))
53         a=cl.Counter(init)
54         l1=a.keys()
55         l2=a.values()
56
57         return l1, l2
58         # END EDITING HERE
59
60     def get_reward(self, arm_rewards):
61         # START EDITING HERE
62         for arm_index,rewards in arm_rewards.items():
63             a = np.sum(rewards)
64             b = len(rewards)
65             self.succs[arm_index]+=a
66             self.fail[arm_index]+=(b-a)
67
68         # END EDITING HERE
```

This is the code of batch sampling.

The plot of Regret vs Horizon for batch sampling looks like this



This is the plot of batch sampling.

3 Task3

3.1 Problem Statement

This task involves dealing with a bandit instance where the horizon is equal to the number of arms. So, for example, if there are 100 arms, then you are only allowed to pull 100 times. However, you are given that the arm means are distributed regularly (in arithmetic progression) between 0 and $(1 - 1/\text{numArms})$.

You need to come up with an algorithm to handle this situation effectively: can you do better than sampling each arm once? Implement your algorithm by editing the `task3.py` file. The APIs you need to implement are essentially the same as `task1.py`.

Once again, you can use `simulator.py` to see regrets as a function of horizon. You may also run `autograder.py` to evaluate your algorithm. Note that even if your algorithm passes the autograder tests, it might fail on the undisclosed tests that are used for evaluation. So you must not hardcode your method to make it work for only the given test instances. For this task, you will again plot regret against horizon (which is the same as the number of arms).

3.2 Many Arms

The problem where the number of arms and horizon equal is quite an interesting one. At first it might be tempting to apply UCB or Thompson or any other sub linear regret algorithms. But those are efficient only when the horizon is much greater than number of arms. So in this case it is best to use the basic epsilon greedy method. In here we explore for some amount of time and based on our exploration we exploit the remaining time. let us see the code of it.

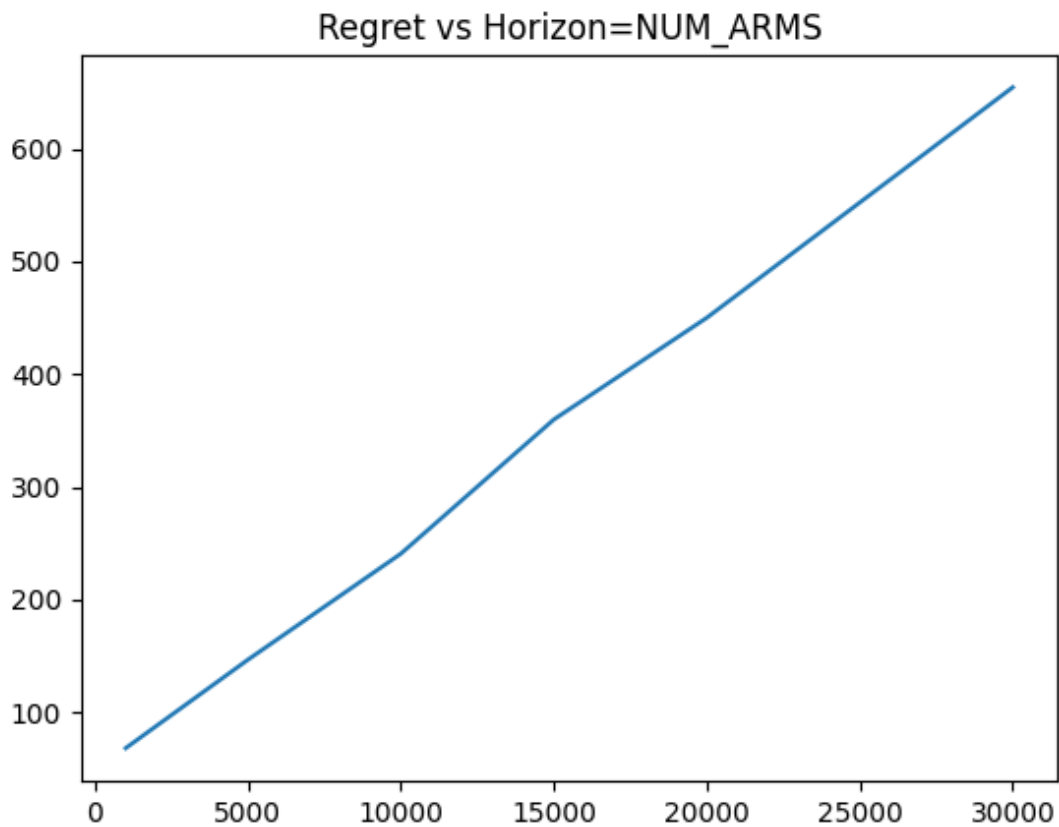
```

28 class AlgorithmManyArms:
29     def __init__(self, num_arms, horizon):
30         self.num_arms = num_arms
31         # Horizon is same as number of arms
32         self.eps = 0.01
33         self.hz = self.eps*horizon
34         self.pulls=0
35         self.counts = np.zeros(num_arms)
36         self.values = np.zeros(num_arms)
37
38     def give_pull(self):
39         # START EDITING HERE
40         if self.pulls < self.hz:
41             return np.random.randint(self.num_arms)
42         else:
43             return np.argmax(self.values)
44
45         #raise NotImplementedError
46         # END EDITING HERE
47
48     def get_reward(self, arm_index, reward):
49         self.pulls+=1
50         self.counts[arm_index] += 1
51         #n = self.counts[arm_index]
52         #value = self.values[arm_index]
53         #new_value = ((n - 1) / n) * value + (1 / n) * reward
54         self.values[arm_index] = ((self.counts[arm_index]-1)/self.counts[arm_index])*self.values[arm_index] +(1/self.counts[arm_index])*reward
55         # START EDITING HERE
56         #raise NotImplementedError
57

```

This is the code of ManyArms.

The plot of Regret vs Horizon for ManyArms looks like this



This is the plot of ManyArms.