

CS 747 : Assignment 1

Report

Roll No. 15D170013

Notes and answers to question :

1) Observations about the results:

In general we see the trend that the

Regret for $\epsilon\text{-greedy} > \text{UCB} > \text{KLUCB} > \text{Thompson Sampling}$

In the $\epsilon\text{-greedy}$ methods : increase in $\epsilon \rightarrow$ increase in regret

2) Patterns Displayed:

The graphs are increasing as expected and the ordering on an average as mentioned above. The 10,000 ones show varied trend after time of around 4000.

3) Effect of epsilon and intuition:

Increasing ϵ increases the regret. This because exploration is higher and for greater values of Horizon (1000, 10000) exploitation seems like a better way forward.

4) Unexpected Results:

As can be seen from the graphs increasing the ϵ results in increase in the regret. The Thompson sampling as read online has lower regrets practically but in my case the KLUCB algorithm performs better in some instances.

Sometimes for smaller regrets and horizons the $\epsilon\text{-greedy}$ method performs better than the remaining but on an average for different horizons. Also in the histogram plots with a horizon of 10,000 we see that Thompson's regret is higher for time < 4000 but reverses and beats KLUCB and $\epsilon\text{-greedy}$ for higher regrets.

5) Further improve performance?

Optimizations : For the ϵ we should have some way in the exploitation to sample certain random arms initially. Also in the client side if we could have implemented horizon dependent ϵ s for better exploitation vs exploration. Also parallelising the code to run on different ports would result in faster improvements.

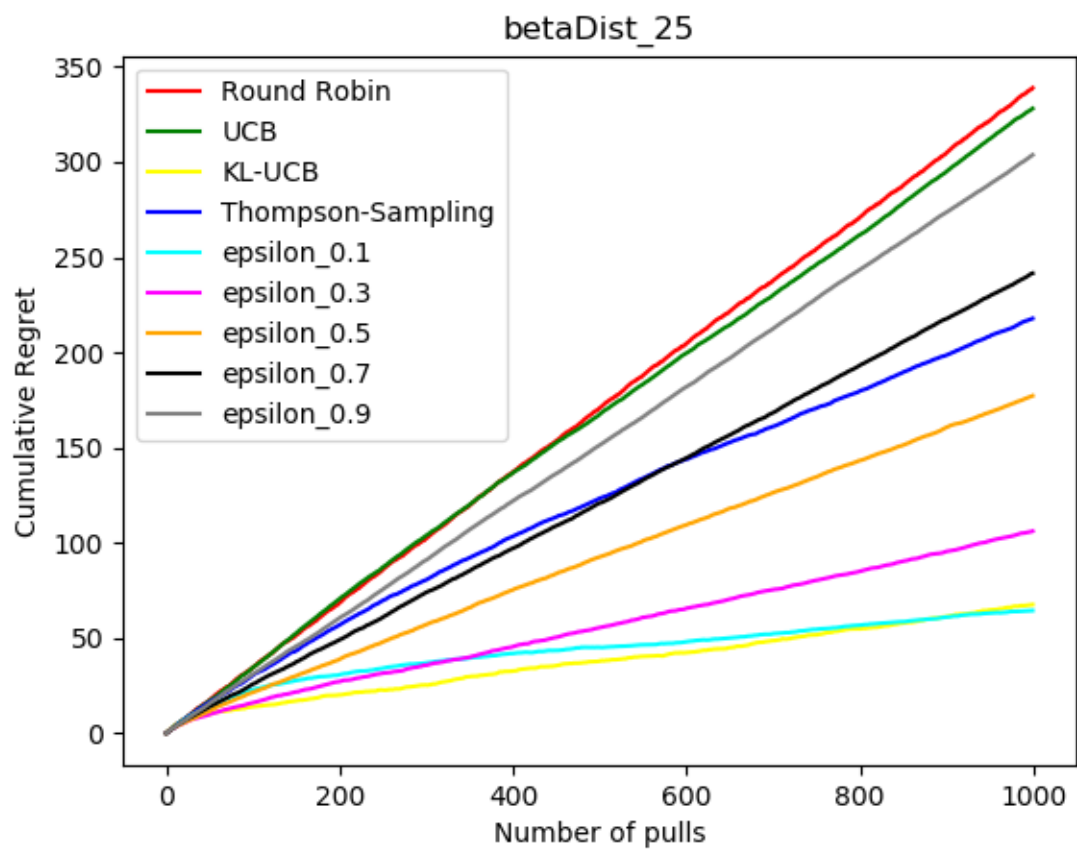
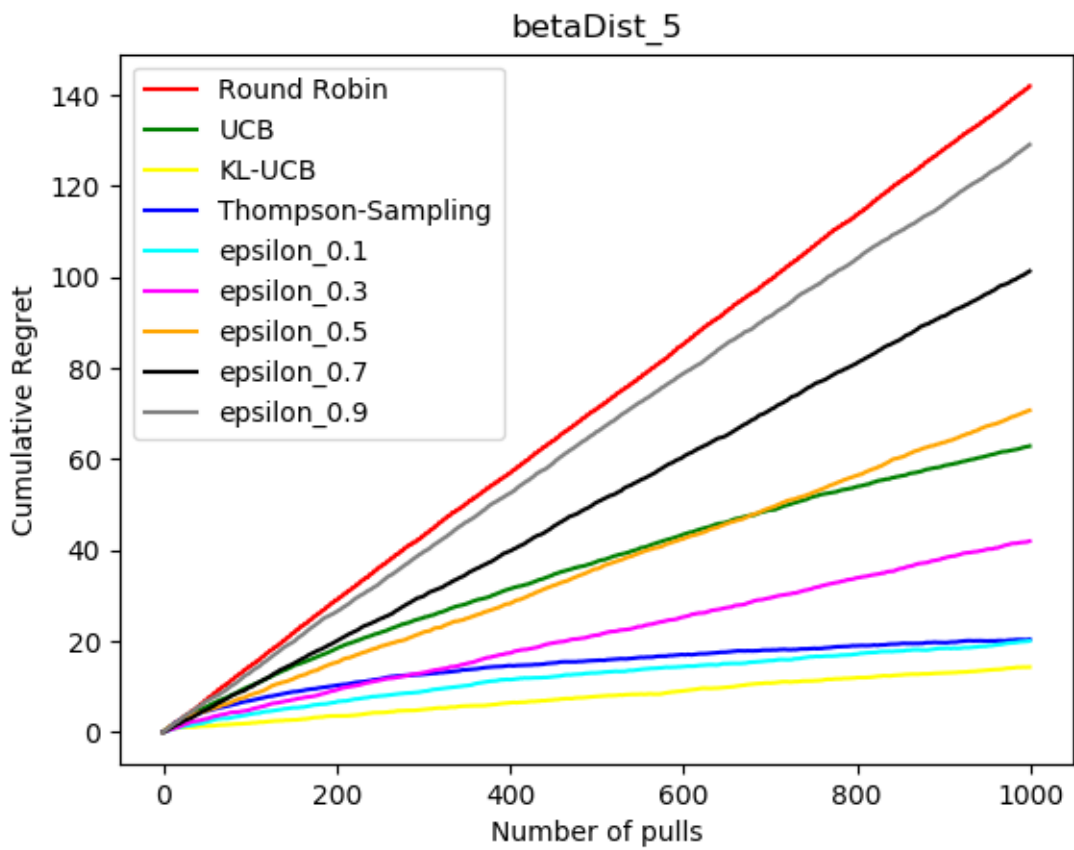
6) Changes for Non Bernoulli reward :

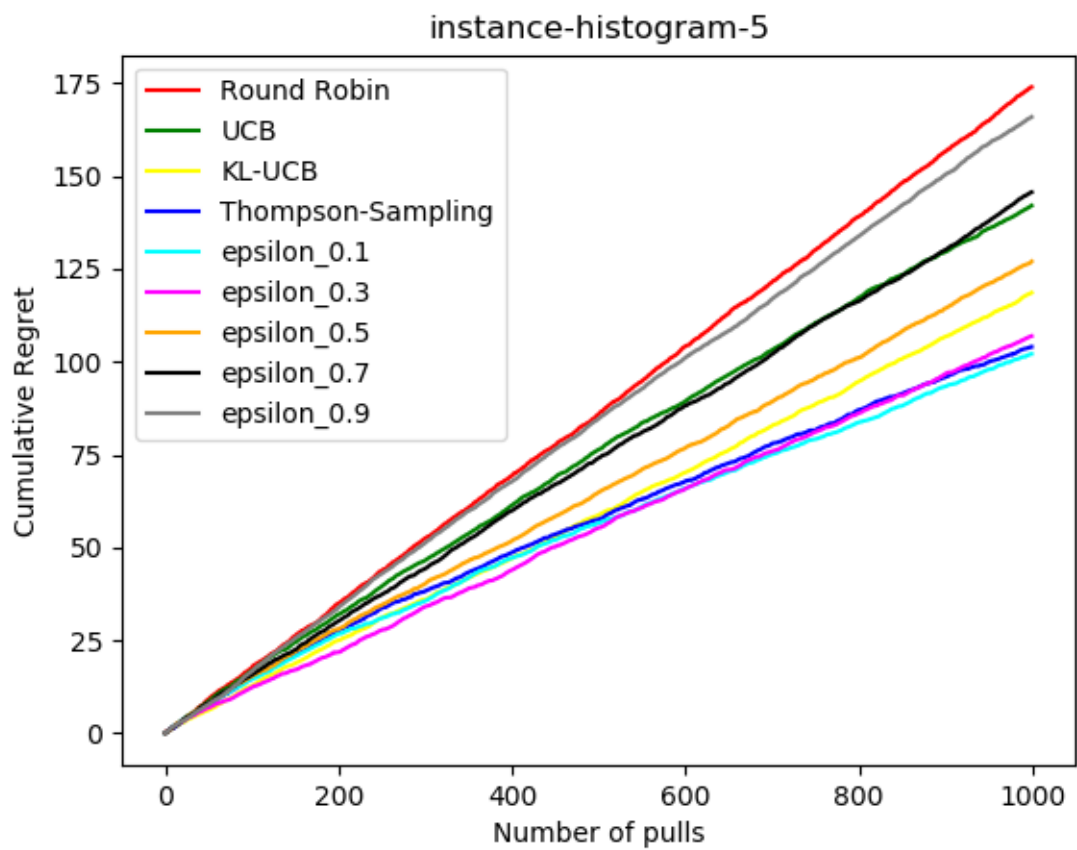
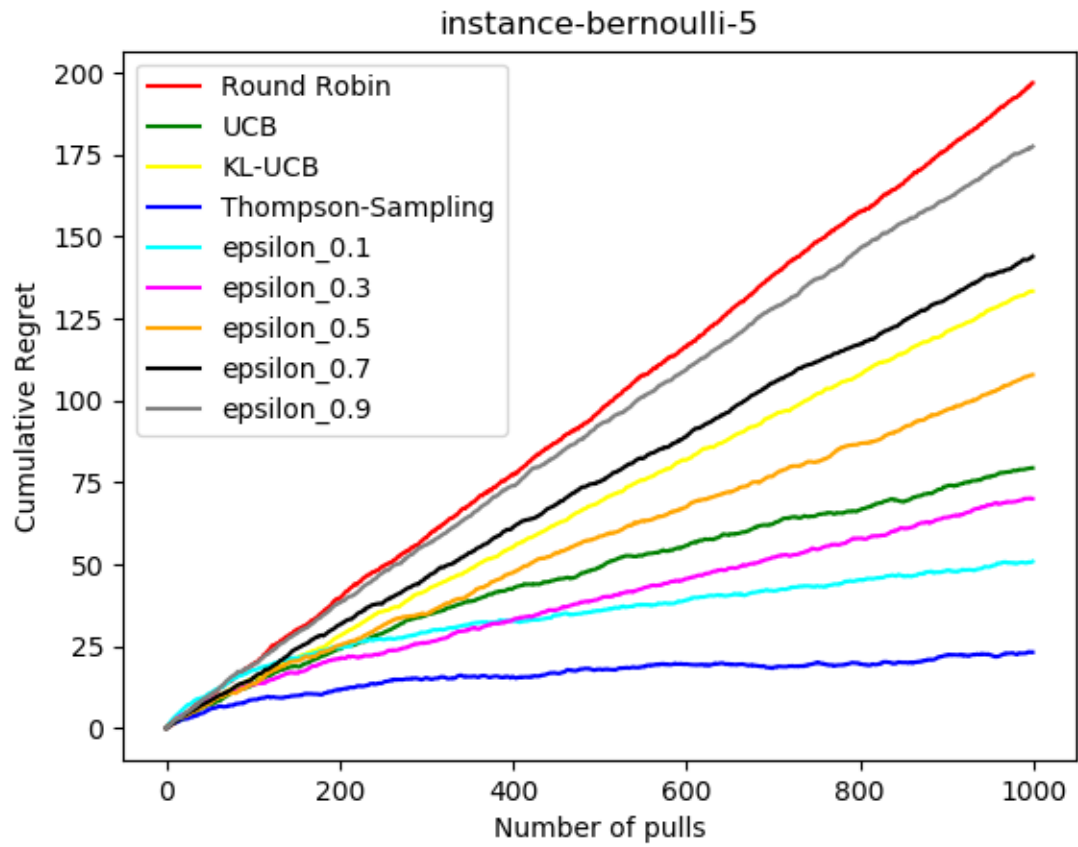
Thomson's Sampling : Handling the cases of non Bernoulli rewards is done in the following way. We sample from a Bernoulli distribution keeping the reward received as the parameter. If the sample is 1 then we consider it as a success and a failure vice versa. For the other algorithms such cases are inherently handled.

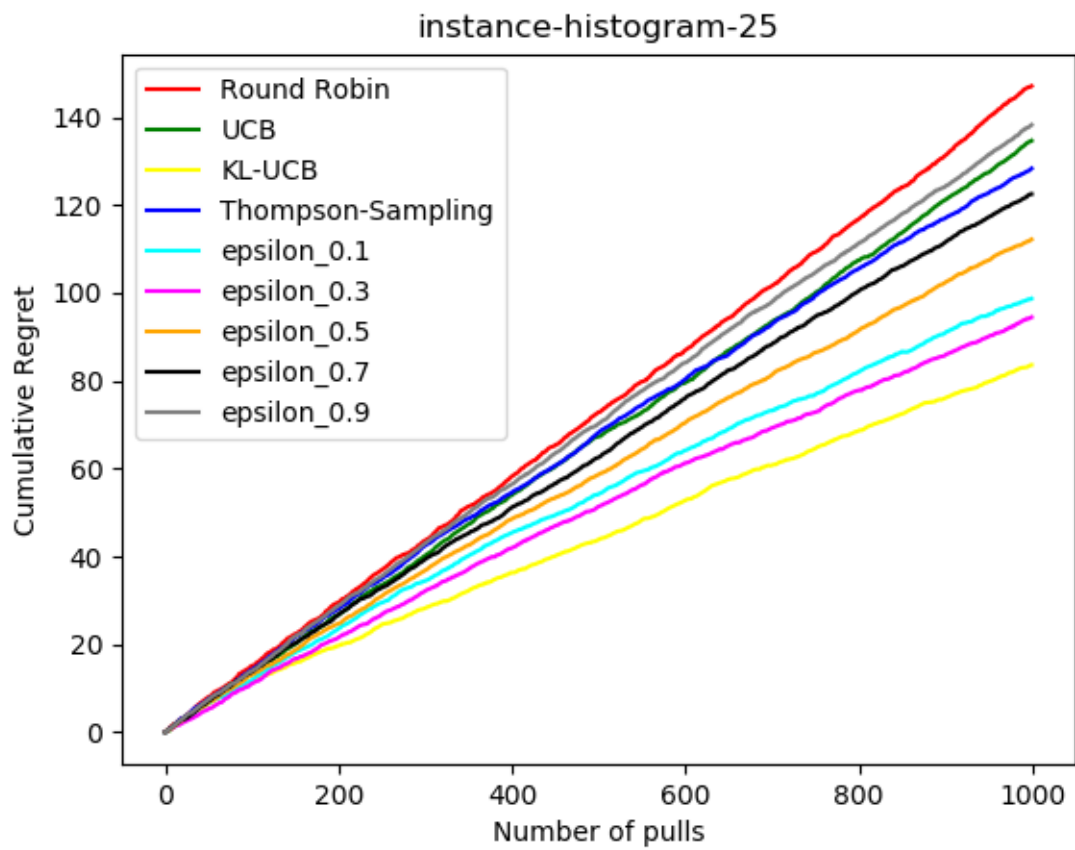
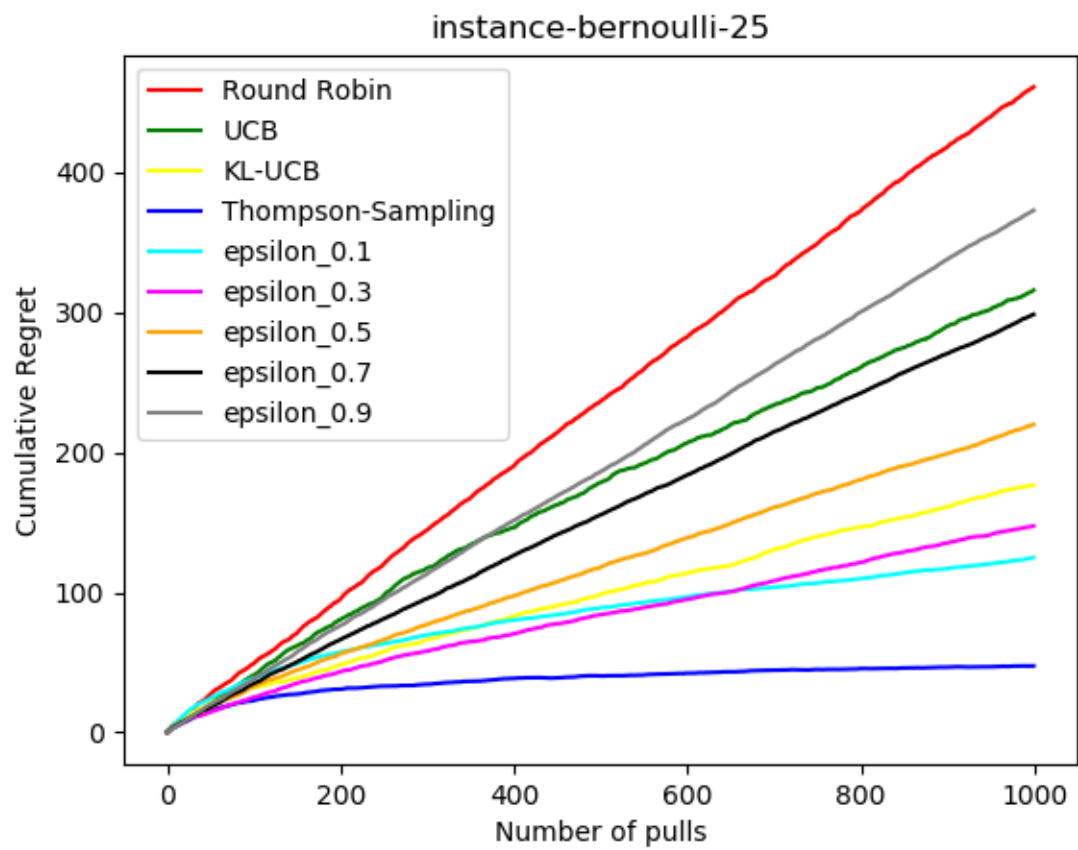
7) Implementation Details Specific:

KLUCB : For running the algo the step of finding q is done using the Binary Search between \hat{p} and 1 and not using Newton iteration method as suggested in the paper. This is done up to an accuracy of 10^{-6} . i.e. 20 loops.

Following are the graphs for horizon set to 1000:







Following are the graphs for horizon set to 10000:

