

# Assignment 1

EE675: Introduction to Reinforcement Learning

January 22, 2025

## Instructions

- Kindly name your submission files as 'RollNo\_Name\_A1\_PartA/B.ipynb', based on the part you are submitting. Marks will be deducted for all submissions that do not follow the naming guidelines.
- You must work out your answers and submit only the iPython Notebook. The code should be well commented and easy to understand as there are marks for this.
- Submissions are to be made through HelloIITK portal. Submissions made through mail will not be graded.
- Answers to the theory questions, if any, should be included in the notebook itself. While using special symbols use the  $\text{\LaTeX}$  mode
- Make sure your plots are clear and have title, legends and clear lines, etc.
- Plagiarism of any form will not be tolerated. If your solutions are found to match with other students or from other uncited sources, there will be heavy penalties and the incident will be reported to the disciplinary authorities.
- In case you have any doubts, please contact these TAs for help: Bhavaj (bhavajs21@iitk.ac.in) or Swastik (swastiks21@iitk.ac.in)

## Part-A

**(Chernoff Bound)** [15 Marks] Suppose  $X_1, X_2, \dots, X_n$  are i.i.d. copies of a  $\mathcal{N}(0, \sigma^2)$  r.v. Then for  $X = \frac{1}{n} \sum_{i=1}^n X_i$  we know that

$$P[X \geq \varepsilon] \leq \exp\left(\frac{-n\varepsilon^2}{2\sigma^2}\right)$$

Write a Python code to run Monte Carlo simulations that verify the inequality. Specifically, for a given  $\varepsilon$  and  $\sigma$ , generate  $n$  samples from the zero mean Gaussian distribution  $x_1, x_2, \dots, x_n$  and check whether the sample average is more than  $\varepsilon$ . Repeat this experiment 500 times and observe in how many experiments out of those 500 experiments, the sample average is more than  $\varepsilon$ . This will give us an empirical estimate of  $P[X \geq \varepsilon]$ .

1. Take  $\sigma = 10$ ,  $\varepsilon = 0.05$ , and plot the empirical estimate as a function of  $n \in [1, 100]$ . In the same plot, include the Chernoff upper bound as a function of  $n$ .
2. Repeat the experiment with the Chebyshev inequality bound and compare it with the Chernoff bound.
3. Repeat the experiment with other values of  $\sigma$  and  $\varepsilon$  and comment on what are your observations.

## Part-B

**(Multi arm bandits | Explore-then-Commit and UCB)** [25 Marks] Consider a two-armed Bernoulli bandit scenario with true means given by  $\mu_1 = \frac{1}{2}$ ,  $\mu_2 = \frac{1}{2} + \Delta$ , for some  $\Delta < \frac{1}{2}$ . Let the time horizon be  $T = 10000$ .

1. Take  $\Delta = \frac{1}{4}$  and run the Monte Carlo simulations to estimate the expected regret of the ETC algorithm which explores each arm  $m = T^{2/3}(\log T)^{1/3}$  times before committing. Specifically, you run the ETC algorithm to compute the sample regret

$$\mu_2 \cdot T - \sum_{t=1}^T \mu(a_t),$$

where  $a_t$  is the arm played in time step  $t$ .

Repeat this experiment 500 times and estimate the expected regret by taking the average of the sample regrets you obtained in all those 500 experiments. [5 Marks]

2. Repeat the above for various values of  $\Delta \in \{0.05, 0.1, 0.2, 0.3, 0.4, 0.45\}$  and plot the estimated regret as a function of  $\Delta$  and verify whether it satisfies the regret upper bound we derived in class. [5 Marks]
3. Repeat the experiment with the UCB algorithm and plot the comparison with ETC. [10 Marks]
4. In the ETC algorithm, assume that we know  $\Delta$ . Fix an  $m$  such that the number of exploration samples is sufficient to make  $\varepsilon < \frac{\Delta}{2}$  with a high probability of  $1 - \frac{1}{T}$ . Repeat the same experiments as above and compare them with UCB. What did you observe? [5 Marks]