

## Assignment 2: March 19

**Instructions:** Discussion among the class participants is highly encouraged. But please make sure that you understand the algorithms and write your own code.

**Submit Solutions by 11:59PM, 31<sup>st</sup> March on Moodle.** Late submission will not be evaluated.

**Question 1 (5 points)** Suppose that  $X$  is  $\sigma$ -subgaussian and  $X_1$  and  $X_2$  are independent and  $\sigma_1$  and  $\sigma_2$ -subgaussian respectively, then:

1.  $\mathbb{E}[X]=0$  and  $\text{Var}[X] \leq \sigma^2$ .
2.  $cX$  is  $|c|\sigma$ -subgaussian for all  $c \in \mathbb{R}$ .
3.  $X_1 + X_2$  is  $\sqrt{\sigma_1^2 + \sigma_2^2}$ -subgaussian.

**Question 2 (10 points)** Suppose that  $X$  is zero-mean and  $X \in [a, b]$  almost surely for constants  $a < b$ .

1. Show that  $X$  is  $(b-a)/2$ -subgaussian.
2. Using Cramer-Chernoff method shows that if  $X_1, X_2, \dots, X_n$  are independent and  $X_t \in [a_t, b_t]$  almost surely with  $a_t < b_t$  for all  $t$ , then prove

$$\mathbb{P}\left(\sum_{t=1}^n (X_t - \mathbb{E}[X_t]) \geq \epsilon\right) \leq \exp\left(-\frac{2\epsilon^2}{\sum_{t=1}^n (b_t - a_t)^2}\right)$$

**Question 3 (5 points)** [Expectation of maximum] Let  $X_1, \dots, X_n$  be a sequence of  $\sigma$ -subgaussian random variables (possibly dependent) and  $Z = \max_{t \in [n]} X_t$ . Prove that

1.  $\mathbb{E}[Z] \leq \sqrt{2\sigma^2 \log(n)}$ .
2.  $\mathbb{P}(Z \geq \sqrt{2\sigma^2 \log(n/\delta)}) \leq \delta$  for any  $\delta \in (0, 1)$ .

**Question 4 (20 points)** [Bernstein's inequality] Let  $X_1, \dots, X_n$  be a sequence of independent random variables with  $X_t - \mathbb{E}[X_t] \leq b$  almost surely and  $S = \sum_{t=1}^n (X_t - \mathbb{E}[X_t])$  and  $v = \sum_{t=1}^n \text{Var}[X_t]$ .

1. Show that  $g(x) = \frac{1}{2} + \frac{x}{3!} + \frac{x^2}{4!} + \dots = \frac{(\exp(x)-1-x)}{x^2}$  is increasing.
2. Let  $X$  be a random variable with  $\mathbb{E}[X] = 0$  and  $X \leq b$  almost surely. Show that  $\mathbb{E}[\exp(X)] \leq 1 + g(b)V[X]$ .
3. Prove that  $(1+\alpha) \log(1+\alpha) - \alpha \geq \frac{3\alpha^2}{6+2\alpha}$  for all  $\alpha \geq 0$ . Prove that this is the best possible approximation in the sense that the 2 in the denominator cannot be increased.

4. Let  $\epsilon > 0$  and  $\alpha = b\epsilon/v$  and prove that

$$\begin{aligned}\mathbb{P}(S \geq \epsilon) &\leq \exp\left(-\frac{v}{b^2}((1+\alpha)\log(1+\alpha) - \alpha)\right) \\ &\leq \exp\left(-\frac{\epsilon^2}{2v\left(1 + \frac{b\epsilon}{3v}\right)}\right).\end{aligned}$$

5. Use the previous result to show that

$$\mathbb{P}\left(S \geq \sqrt{2v \log\left(\frac{1}{\delta}\right)} + \frac{2b}{3} \log\left(\frac{1}{\delta}\right)\right) \leq \delta.$$

**Question 5 (10 points)** Show that

$$R_T \leq \min\left\{T\Delta, \Delta + \frac{4}{\Delta} \left(1 + \max\left\{0, \log\left(\frac{T\Delta^2}{4}\right)\right\}\right)\right\}$$

implies the regret of an optimally tuned Explore-then-Commit (ETC) algorithm for subgaussian 2-armed bandits with means  $\mu_1, \mu_2 \in \mathbb{R}$  and  $\Delta = |\mu_1 - \mu_2|$ , satisfies  $R_T \leq \Delta + C\sqrt{T}$  where  $C > 0$  is a universal constant.

**Question 6 (5 points)** Fix  $\delta \in (0, 1)$ . Modify the ETC algorithm to depend on  $\delta$  and prove a bound on the pseudo-regret  $R_T = T\mu^* - \sum_{t=1}^T u_{A_t}$  of ETC algorithm that holds with probability  $1 - \delta$  where  $A_t$  is the arm chosen in the round  $t$ .

**Hint:** Choose ‘ $m$ ’ appropriately in the regret upper bound of ETC algorithm which is proved in the class.

**Question 7 (5 points)** Fix  $\delta \in (0, 1)$ . Prove a bound on the random regret  $R_T = T\mu^* - \sum_{t=1}^T X_t$  of ETC algorithm that holds with probability  $1 - \delta$ . Compare this to the bound derived for the pseudo-regret in the question 5. What can you conclude?

**Question 8 (10 points)** Assume the rewards are 1-subgaussian and there are  $k \geq 2$  arms. The  $\epsilon$ -greedy algorithm depends on a sequence of parameters  $\epsilon_1, \epsilon_2, \dots$ . First it chooses each arm once and subsequently chooses  $A_t = \arg \max_i \hat{\mu}_i(t-1)$  with probability  $1 - \epsilon_t$  and otherwise chooses an arm uniformly at random.

1. Prove that if  $\epsilon_t = \epsilon > 0$ , then  $\lim_{T \rightarrow \infty} \frac{R_T}{T} = \frac{\epsilon}{k} \sum_{i=1}^k \Delta_i$ .

2. Let  $\Delta_{\min} = \min\{\Delta_i : \Delta_i > 0\}$  where  $\Delta_i = \mu^* - \mu_i$ , and  $\epsilon_t = \min\left\{1, \frac{Ck}{t\Delta_{\min}^2}\right\}$  where  $C > 0$  is a sufficiently large universal constant. Prove that there exists a universal  $C' > 0$  such that

$$R_T \leq C' \sum_{i=1}^k \left( \Delta_i + \frac{\Delta_i}{\Delta_{\min}^2} \log \max\left\{e, \frac{T\Delta_{\min}^2}{k}\right\} \right).$$

**Question 9 (10 points)** Fix a 1-subgaussian  $k$ -armed bandit environment and a horizon  $T$ . Consider the version of UCB that works in phases of exponentially increasing length of  $1, 2, 4, \dots$ . In each phase, the algorithm uses the action that would have been chosen by UCB at the beginning of the phase.

1. State and prove a bound on the regret for this version of UCB.
2. How would the result change if the  $l^{\text{th}}$  phase had a length of  $\lceil \alpha^l \rceil$  with  $\alpha > 1$ ?

**Submission Format and Evaluation:** You should submit a report along with your code. Please zip all your files and upload via Moodle. The zipped folder should be named as YourRegistrationNo.zip e.g. '154290002.zip'. The report should contain one figure with four plots corresponding to each algorithm in Q.1. Write a brief summary of your observations. We may also call you to a face-to-face session to explain your code.