

### Exercise 3: The $(\alpha, \psi)$ -UCB algorithm (can be solved after Course #4)

Let  $\psi : \mathbb{R} \rightarrow \mathbb{R}$  be a convex function such that  $\psi(x) = \psi(-x)$  for all  $x \in \mathbb{R}$ . Consider a bandit model  $\mathcal{D}$  such that for all  $\nu \in \mathcal{D}$ , if  $X$  denotes a random variable with distribution  $\nu$ , then

$$\forall \lambda \geq 0, \quad \max \left\{ \ln \mathbb{E}_\nu \left[ e^{\lambda(X - \mathbb{E}[X])} \right], \ln \mathbb{E}_\nu \left[ e^{\lambda(\mathbb{E}[X] - X)} \right] \right\} \leq \psi(\lambda). \quad (\star)$$

For all  $x \geq 0$ , we define the convex conjugate of  $\psi$ ,

$$\psi^*(x) = \sup \{ \lambda x - \psi(\lambda) : \lambda \geq 0 \},$$

and assume that  $\psi^*$  is invertible, with inverse denoted by  $(\psi^*)^{-1}$ .

1. Provide such a function  $\psi$  for the model  $\mathcal{D} = \mathcal{P}([0, 1])$  of all probability distributions over  $[0, 1]$ . Compute  $\psi^*$  and its inverse.

We generalize the UCB algorithm for stochastic bandits in the following way. We consider the same setting and use the same notation as the ones used in class and in Exercise 3 of the present statement: a stochastic bandit problem is formed by  $K \geq 2$  probability distributions  $\nu_1, \dots, \nu_K$  in  $\mathcal{D}$  with respective expectations  $\mu_k$ , their maximal expectation is denoted by  $\mu^*$ , the gap of arm  $k$  is  $\Delta_k = \mu^* - \mu_k$ , etc.

---

#### $(\alpha, \psi)$ -UCB algorithm

---

*Parameters:*  $\alpha > 0$  and  $\psi : \mathbb{R} \rightarrow \mathbb{R}$  with  $\psi(x) = \psi(-x)$  for all  $x \geq 0$

*Initialization:* Play each arm once, i.e.,  $I_t = t$  for  $t \in \{1, \dots, K\}$ , get a reward  $Y_t \sim \nu_t$

For  $t \geq K + 1$ ,

1. Compute, for all  $k \in \{1, \dots, K\}$ ,

$$N_k(t-1) = \sum_{s=1}^{t-1} \mathbb{1}_{\{I_s=k\}} \quad \text{and} \quad \hat{\mu}_{k,t-1} = \frac{1}{N_k(t-1)} \sum_{s=1}^{t-1} Y_s \mathbb{1}_{\{I_s=k\}}$$

2. Pick an arm (ties broken arbitrarily)

$$I_t \in \operatorname{argmax}_{k \in \{1, \dots, K\}} \left\{ \hat{\mu}_{k,t-1} + (\psi^*)^{-1} \left( \frac{\alpha \ln t}{N_k(t-1)} \right) \right\}$$

3. Get a reward  $Y_t \sim \nu_{I_t}$  (conditionally to  $I_t$ )
- 

We want to upper bound the pseudo-regret of the  $(\alpha, \psi)$ -UCB algorithm as follows: for  $\alpha > 2$ ,

$$\bar{R}_T = T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T Y_t \right] \leq \sum_{k: \Delta_k > 0} \Delta_k \left( \frac{\alpha}{\psi^*(\Delta_k/2)} \ln T + \frac{2\alpha}{\alpha - 2} \right). \quad (\text{B})$$

To that end, we first show that for each arm  $k$  and  $t \geq K + 1$ , an upper confidence bound on  $\mu_k$  is given by

$$\hat{\mu}_{k,t-1} + (\psi^*)^{-1} \left( \frac{\alpha \ln t}{N_k(t-1)} \right).$$

2. Prove that for all  $t \geq 1$  and all  $\lambda \geq 0$ ,

$$\mathbb{E} \left[ \exp \left( -\lambda (Y_t - \mu_k) \mathbb{1}_{\{I_t=k\}} \right) \middle| \mathcal{F}_{t-1} \right] \leq \exp(\psi(\lambda) \mathbb{1}_{\{I_t=k\}})$$

for a filtration  $\mathcal{F} = (\mathcal{F}_t)_{t \geq 0}$  to specify explicitly.

Construct an  $\mathcal{F}$ -adapted supermartingale  $(M_t)_{t \geq 0}$  based on this inequality.

3. Prove that for all  $t \geq K + 1$ , all  $\ell \geq 1$ , and all  $\varepsilon > 0$ ,

$$\mathbb{P}\left\{\widehat{\mu}_{k,t-1} + \varepsilon \leq \mu_k \quad \text{and} \quad N_k(t-1) = \ell\right\} \leq \exp(-\ell \psi^*(\varepsilon)).$$

4. Provide a bound, for  $t \geq K + 1$ , on

$$\mathbb{P}\left\{\widehat{\mu}_{k,t-1} + (\psi^*)^{-1}\left(\frac{\alpha \ln t}{N_k(t-1)}\right) \leq \mu_k\right\}.$$

5. Briefly indicate how to bound, for  $t \geq K + 1$ ,

$$\mathbb{P}\left\{\widehat{\mu}_{k,t-1} - (\psi^*)^{-1}\left(\frac{\alpha \ln t}{N_k(t-1)}\right) > \mu_k\right\}.$$

To establish the regret bound, we first fix a suboptimal arm  $j$  and an optimal arm  $a^*$ .

6. Explain why  $I_t = j$  for  $t \geq K + 1$  entails one of the following events:

$$\begin{aligned} & \widehat{\mu}_{a^*,t-1} + (\psi^*)^{-1}\left(\frac{\alpha \ln t}{N_{a^*}(t-1)}\right) \leq \mu^*, \\ \text{or} & \quad \widehat{\mu}_{j,t-1} - (\psi^*)^{-1}\left(\frac{\alpha \ln t}{N_j(t-1)}\right) > \mu_j, \\ \text{or} & \quad N_j(t-1) < \frac{\alpha \ln t}{\psi^*(\Delta_j/2)}. \end{aligned}$$

7. Establish the regret bound (B).

We conclude this exercise with a discussion of the bound for the model  $\mathcal{D} = \mathcal{P}([0, 1])$ .

8. Provide also a distribution-free bound for  $(\alpha, \psi)$ -UCB on this model, i.e., a bound over all distributions satisfying  $(\star)$ . You need first to think of a suitable value for  $\alpha$ .