

---

## Supplementary Material: Optimal Algorithms for Lipschitz Bandits with Heavy-tailed Rewards

---

### A. Proof of Theorem 1

**Theorem 1** Assume (1) and (2) hold. For sufficiently large  $T$  such that

$$\log T \geq \frac{5}{8}(4r)^{-\frac{1}{\epsilon}}$$

the regret of SDTM with parameter  $r > 0$  satisfies

$$\mathbb{E}[R(T)] \leq 2rT + (4rT)^{\frac{1}{1+\epsilon}} (16N_c(r) \log T)^{\frac{\epsilon}{1+\epsilon}}$$

where  $N_c(r)$  is the  $r$ -covering number of the arm set  $\mathcal{X}$ .

**Proof.** Let  $x_* \in \arg \max_{x \in \mathcal{X}} \mu(x)$  be an optimal arm. By the definition of the oracle, there must exist  $k \in [K]$  such that  $x_* \in \mathcal{X}_k$  and hence  $D(x_*, \bar{X}_k) \leq 2r$ . Since the expected reward function is Lipschitz, we have

$$\mu(x_*) - \mu(\bar{X}_k) \leq D(x_*, \bar{X}_k) \leq 2r. \quad (17)$$

On the other hand, let  $\bar{X}_* \in \arg \max_{x_i, i \in [K]} \mu(x_i)$  be an optimal skeleton arm. By theoretical guarantees of UCB policies used with the truncated mean estimator (Bubeck et al. 2013, Proposition 1), the expected difference between the cumulative reward of the pulled arms and that of the optimal skeleton arm  $\bar{X}_*$  can be upper bounded as follows

$$\mathbb{E} \sum_{t=1}^T \mu(x_*) - \sum_{t=1}^T \mu(x_t) \leq (4rT)^{\frac{1}{1+\epsilon}} (16K \log T)^{\frac{\epsilon}{1+\epsilon}}. \quad (18)$$

Combining (17) and (18) and recalling that  $K \leq N_c(r)$ , we obtain

$$\mathbb{E} \sum_{t=1}^T \mu(x_*) - \sum_{t=1}^T \mu(x_t) \leq 2rT + (4rT)^{\frac{1}{1+\epsilon}} (16N_c(r) \log T)^{\frac{\epsilon}{1+\epsilon}}$$

where we use the fact that  $\mu(\bar{X}_k) \geq \mu(x_*)$ .

### B. Proof of Corollary 2

**Corollary 2** We have

$$\sum_{r=2^{-i}: i \in \mathbb{N}, r \geq r_0} \frac{N_z(r)}{r^{1/\epsilon}} = O(r_0^{-(d_z+1/\epsilon)})$$

and thus

$$R(T) = O\left(\inf_{r_0 \in (0,1)} r_0 T + \log T \cdot r_0^{-(d_z+1/\epsilon)}\right) \\ \leq T^{\frac{d_z \epsilon + 1}{d_z \epsilon + \epsilon + 1}}$$

where  $d_z$  is the zooming dimension of  $(\mathcal{X}; \mu)$ , defined in (5).

**Proof.** We have

$$\sum_{r=2^{-i}: i \in \mathbb{N}, r \geq r_0} \frac{N_z(r)}{r^{1/\epsilon}} = \sum_{i \in \mathbb{N}: 2^{-i} \geq r_0} 2^{id_z + i/\epsilon} Z \sum_{i=0}^{\lfloor \log_2 \frac{1}{r_0} \rfloor} 2^{id_z + i/\epsilon} Z \sum_{i=0}^{\lfloor \log_2 \frac{1}{r_0} \rfloor + 1} 2^{id_z + i/\epsilon} Z di \frac{(2^{-i} r_0)^{d_z + 1/\epsilon} Z}{\log(2^{d_z + 1/\epsilon})}$$

where  $Z$  is the zooming constant of  $(\mathcal{X}; D)$ .

### C. Proof of Theorem 3

**Theorem 3** Assume (2) and (12) hold. With probability at least  $1 - 2^{-\epsilon}$ , the regret of ADMM satisfies

$$R(T) \leq \inf_{r_0 \in (0,1)} r_0 T + 68(102^{-\epsilon})^{\frac{1}{\epsilon}} \log(e^{1/8} T^{2-\epsilon}) \times \max_{r=2^{-i}; i \in \mathbb{N}, r \geq r_0} \frac{N_z(r)}{r^{1/\epsilon}}$$

where  $\bar{\cdot}$  is defined in (13) and  $N_z(r)$  is the  $r$ -zooming number of  $(X; \bar{\cdot})$ . Furthermore, by the first inequality in Corollary 2 we have

$$R(T) \leq \Theta \left( T^{\frac{d_z \epsilon + 1}{d_z \epsilon + \epsilon + 1}} \right)$$

where  $d_z$  is the zooming dimension of  $(X; \bar{\cdot})$ , defined in (5).

**Proof.** We use the same notations as in the proof of Theorem 2 and propose the following lemmas, which are counterparts of Lemmas 1, 2, 3, and 4 respectively. For brevity, we only prove Lemmas 5 and 6, and the proofs of Lemmas 7 and 8 can be done in the same way as in appendices G and H respectively.

**Lemma 5** Let  $R$  be the set comprised of all rounds in which Step 21 of Algorithm 4 is executed. Then, with probability at least  $1 - 2^{-\epsilon}$ , for all rounds  $t \in R$  and all active arms  $x \in A_t$ , we have

$$j_t(x) \leq (x) j_{t+1}(x):$$

**Proof.** Fix  $t \in R$  and  $x \in A_t$ . By Lemma 2 in Bubeck et al. (2013), with probability at least  $1 - \frac{2\delta}{T^2}$  we have

$$j_t(x) \leq (x) j_{t+1}(x) \frac{(12^{-\epsilon})^{\frac{1}{1+\epsilon}}}{n_t(x)} \frac{16 \log(e^{1/8} T^{2-\epsilon})^{\frac{\epsilon}{1+\epsilon}}}{n_t(x)} \frac{(12^{-\epsilon})^{\frac{1}{1+\epsilon}}}{n_t(x)} \frac{16 \log(e^{1/8} T^{2-\epsilon})^{\frac{\epsilon}{1+\epsilon}}}{n_t(x)} = r_{t+1}(x):$$

Taking the union bound over  $x \in A_t$  and  $t \in R$  and noticing  $|A_t| \leq T; \delta \in R$ , we conclude the proof.

**Lemma 6** With probability at least  $1 - 2^{-\epsilon}$ , for all rounds  $t \in [T]$  and all active arms  $x \in A_t$ , we have

$$\Delta(x) \leq 3^{\rho_-} 2 r_{t+1}(x):$$

**Proof.** Fix  $t \in [T]$ . For each active arm  $x \in A_t$ , there exist three different scenarios as follows.

(i)  $x$  is pulled by Step 4 or Step 10 of Algorithm 4 in round  $t$ . In this scenario, on one hand, we have

$$n_t(x) \leq 16 \log(e^{1/8} T^{2-\epsilon}) + 1$$

and hence

$$\begin{aligned} r_{t+1}(x) &= (12^{-\epsilon})^{\frac{1}{1+\epsilon}} \frac{16 \log(e^{1/8} T^{2-\epsilon})^{\frac{\epsilon}{1+\epsilon}}}{n_t(x)} \frac{(12^{-\epsilon})^{\frac{1}{1+\epsilon}}}{16 \log(e^{1/8} T^{2-\epsilon}) + 1} \\ &\leq (3^{\rho_-} 2)^{-\frac{1}{1+\epsilon}} \frac{35}{36} \frac{16 \log(e^{1/8} T^{2-\epsilon})^{\frac{\epsilon}{1+\epsilon}}}{36} = \frac{35}{36} 3^{\rho_-} \frac{35}{36} \frac{16 \log(e^{1/8} T^{2-\epsilon})^{\frac{\epsilon}{1+\epsilon}}}{36} \frac{1}{3^{\rho_-} 2} \end{aligned}$$

where the second inequality follows from the definition of  $\bar{\cdot}$  in (13) and the following fact:  $16 \log(e^{1/8} T^{2-\epsilon}) > 35$  for  $T > 1$  and  $\epsilon \in (0; 1/2)$ . On the other hand, let  $x_* \in \arg \max_{x \in \mathcal{X}} (x)$  be an optimal arm. We have  $\Delta(x) = (x_*) - (x) \leq D(x_*; x) \leq 1$ . Thus, we obtain  $3^{\rho_-} 2 r_{t+1}(x) \geq \Delta(x)$ .

(ii)  $x$  is pulled by Step 12 of Algorithm 4 in round  $t$ . In this case, we have  $t - 1 \in R$  and the arm selection rule implies

$$b_{t-1}(x) + 2r_t(x) \leq b_{t-1}(x') + 2r_t(x'); \quad x' \in A_t: \quad (19)$$

Note that  $A_t = A_{t-1}$ . By Lemma 5, we get

$$r_t(x) - b_{t-1}(x) \leq r_t(x') - b_{t-1}(x') \quad (x') \quad r_t(x'); \quad \delta x' \geq A_t: \quad (20)$$

Combining (19) and (20), we obtain

$$r_t(x) + 3r_t(x) \leq r_t(x') + r_t(x'); \quad \delta x' \geq A_t: \quad (21)$$

Note that the execution of Step 12 implies  $\sum_{x \in A_t} B(x; r_t(x))$ . Thus, for the optimal arm  $x_*$ , there must exist an active arm  $\bar{x}_* \in A_t$  such that

$$D(x_*; \bar{x}_*) \leq r_t(\bar{x}_*)$$

which, together with the Lipschitz property of  $r_t$ , indicates

$$r_t(x_*) \leq r_t(\bar{x}_*) + r_t(\bar{x}_*):$$

Combining the above inequality and (21) with substitution  $x' = \bar{x}_*$ , we obtain

$$r_t(x) + 3r_t(x) \leq r_t(x_*):$$

On the other hand, we have

$$\frac{r_{t+1}(x)}{r_t(x)} = \frac{n_{t-1}(x)}{n_t(x)} \frac{1}{1+\epsilon} = \frac{n_{t-1}(x)}{n_{t-1}(x) + 1} \frac{1}{1+\epsilon} \leq \frac{1}{2}.$$

Therefore, we get  $3 \frac{r_{t+1}(x)}{r_t(x)} \leq 3r_t(x) \leq \Delta(x)$ :

(iii)  $x$  is not played in round  $t$ . In this scenario, let  $s$  be the last round in which  $x$  is pulled. Then, we have  $r_{t+1}(x) = r_{s+1}(x)$  and the proof reduces to (i) or (ii).

**Lemma 7** With probability at least  $1 - 2^{-i}$ , for all  $i = 0, 1, 2, \dots$ ,

$$|\bar{A}_T(i)| \leq N_z(2^{-i});$$

**Lemma 8** With probability at least  $1 - 2^{-i}$ , for all  $i = 0, 1, 2, \dots$ ,

$$\sum_{x \in \bar{A}_T(i)} n_T(x) \Delta(x) \leq 2^{\frac{i+1}{\epsilon}} (51)^{\frac{1}{\epsilon}} 68 \log(e^{1/8} T^2) N_z(2^{-i});$$

The remaining proof is the same as that of Theorem 2 and is omitted here.

## D. Proof of Theorem 4

**Theorem 4** Fix an arm set  $X$  with diameter 1 and a parameter of moment  $\epsilon \in (0, 1]$ . Define  $\beta = \frac{2^{1/\epsilon} \epsilon}{\log 2}$  and

$$R_c(T) = \inf_{r_0 \in (0, 1)} r_0 T + \log T \sum_{r=2^{-i}; i \in \mathbb{N}, r \geq r_0} \frac{N_c(r)}{r^{1/\epsilon}}$$

where  $N_c(r)$  is the  $r$ -covering number of  $X$ . Then, for any  $T > 2$  and any positive number  $R \geq R_c(T)$ , there exists a set  $I$  of problem instances on  $X$  such that

(i) for each problem instance  $I \in I$ , define

$$R_z(T) = \inf_{r_0 \in (0, 1)} r_0 T + \log T \sum_{r=2^{-i}; i \in \mathbb{N}, r \geq r_0} \frac{N_z(r)}{r^{1/\epsilon}}$$

in which  $N_z(r)$  is the  $r$ -zooming number of  $I$ . We have  $R_z(T) \leq 3R = (8 \log T) R$ :

(ii) for any algorithm  $A$ , there exists at least one problem instance  $I \in I$  on which the expected regret of  $A$  satisfies  $\mathbb{E}[R(T)] \geq R = (2560 \log T) R$ :

**Proof.** Our proof is inspired by Slivkins (2014) and makes use of the needle-in-the-haystack technique, which is firstly proposed by Auer et al. (2002b) for analyzing multi-armed bandits and then extended to Lipschitz bandits by Kleinberg et al. (2013).

**Step 1 (Constructing instance set  $I$ )** We begin with the following lemma.

**Lemma 9** Define

$$R'_c(T) = \inf_{r_0 \in (0,1)} r_0 T + \log T \frac{N_c(r_0)}{r_0^{1/\epsilon}} ;$$

Then for any  $T > 2$ ,  $R_c(T) \leq R'_c(T)$ .

**Proof.** Since  $N_c(r)$  is non-increasing with  $r$ , we have

$$R_c(T) = \inf_{r_0 \in (0,1)} r_0 T + \log T \times \prod_{r=2^{-i}: i \in \mathbb{N}, r \geq r_0} \frac{N_c(r)}{r^{1/\epsilon}} \leq \inf_{r_0 \in (0,1)} r_0 T + \log T \times N_c(r_0) \times \prod_{r=2^{-i}: i \in \mathbb{N}, r \geq r_0} r^{-\frac{1}{\epsilon}}$$

in which the last term can be upper bounded as follows

$$\prod_{r=2^{-i}: i \in \mathbb{N}, r \geq r_0} r^{-\frac{1}{\epsilon}} = \prod_{i=0}^{\lfloor \log_2 \frac{1}{r_0} \rfloor} 2^{\frac{i}{\epsilon}} = \prod_{i=0}^{\lfloor \log_2 \frac{1}{r_0} \rfloor} 2^{\frac{i}{\epsilon}} = \frac{(2=r_0)^{1/\epsilon}}{\log(2^{1/\epsilon})} \frac{1}{\log 2} \frac{2^{1/\epsilon}}{r_0^{1/\epsilon}} ;$$

Recalling  $\frac{2^{1/\epsilon}}{\log 2} > 1$ , we conclude the proof.

Fix  $T > 2$  and  $R \leq R_c(T)$ . Let  $r = \frac{R}{2\kappa T(1+\log T)}$  and  $N = \max(2; bTr^{1+1/\epsilon})$ . Based on Lemma 9, we can bound  $r$  and  $N$  as follows.

**Lemma 10** We say a subset  $S \subseteq X$  is an  $r$ -packing of  $X$  if the distance between any two points in  $S$  is at least  $r$ , i.e.,  $\inf_{u,v \in S} D(u;v) \geq r$ . Let  $N_p(r)$  denote the  $r$ -packing number of  $X$ , defined as the maximal number of points in an  $r$ -packing of  $X$ :

$$N_p(r) = \max \{ |S| : S \text{ is an } r\text{-packing of } X \}$$

We have

$$r < 1/2 \text{ and } N \leq N_p(r);$$

**Proof.** Define function  $f(r) = \frac{N_c(r)}{r^{1+1/\epsilon}}$ . Since  $f(1) = 1$ ,  $\lim_{r \rightarrow 0} f(r) = +\infty$ , and  $f(r)$  is decreasing on  $(0;1)$ , there must exist  $b \in (0;1)$  such that  $f(b) = \frac{R}{T} = f(b=2)$ . From the first inequality, we obtain

$$R \leq R_c(T) \leq R'_c(T) \leq bT + \frac{N_c(b)}{b^{1/\epsilon}} \log T \leq bT(1 + \log T)$$

which implies  $r = \frac{\widehat{r}T(1+\log T)}{2\kappa T(1+\log T)} = \frac{\widehat{r}}{2} < \frac{1}{2}$ . From the second inequality, we have

$$Tr^{1+1/\epsilon} \leq T(b=2)^{1+1/\epsilon} \frac{N_c(b=2)}{(b=2)^{1+1/\epsilon}} \leq N_c(r);$$

We conclude the proof by the fact that  $N_c(r) \leq N_p(r)$  (Kleinberg et al., 2013) and  $2 \leq N_p(r)$ .

The above lemma ensures that we can find a set of arms  $U = \{u_1; \dots; u_N\} \subseteq X$  such that  $\inf_{x,y \in U} D(x;y) \geq r$ . Based on  $U$ , we construct a set of problem instances  $I = \{I_1; \dots; I_N\}$ . Let us fix  $i \in [N]$  and describe the construction of  $I_i$ : the expected reward function  $\mu_i$  is defined as

$$\mu_i(x) = \begin{cases} \frac{7r}{8}, & x = u_i \\ \frac{3r}{4}, & x = u_j; j \in [N] \text{ and } j \neq i \\ \max(\frac{r}{2}; \max_{u \in U} \mu_i(u) - D(x;u)); & \text{otherwise} \end{cases} \quad (22)$$

and the reward distributions are defined by

$$\Pr(y_j|x) = p_i(y_j|x) = \begin{cases} \mu_i(x)r^{1/\epsilon}, & y = r^{-1/\epsilon} \\ 1 - \mu_i(x)r^{1/\epsilon}, & y = 0 \end{cases}; \quad (23)$$

One can show that for  $i = 1; 2; \dots; N$ ,  $\mu_i$  is Lipschitz and the  $(1 + \frac{1}{\epsilon})$ -th moment of  $p_i$  is upper bounded by  $7/8$ .

**Step 2 (Proving i)** Let  $I$  be a problem instance in  $\mathcal{I}$ . Recall the definition of  $\rho$ -optimal region:  $X_\rho = \{x \in X : \Delta(x) \leq \rho\}$ . It is clear that for  $\rho = 3r/4$ , we have  $X_\rho = \emptyset$  and thus  $N_z(\rho) = 0$ . It follows that

$$R_z(T) = \inf_{r_0 \in (0,1)} r_0 T + \log T \prod_{\rho=2^{-i}: i \in \mathbb{N}, \rho \geq r_0} \frac{N_z(\rho)}{1/\epsilon} \leq \frac{3}{4} r T + \log T \prod_{\rho=2^{-i}: i \in \mathbb{N}, \rho \geq \frac{3}{4}r} \frac{N_z(\rho)}{1/\epsilon} \leq \frac{3R}{8 \log T}$$

where the last inequality is due to  $r = \frac{R}{2\kappa T(1+\log T)}$ .

**Step 3 (Proving ii)** Following the framework of Kleinberg et al. (2013), we first introduce an auxiliary problem instance  $I_0$  in which the expected reward function  $\rho_0$  is defined as

$$\rho_0(x) = \begin{cases} \frac{3r}{4}; & x = u_j; j \in [N] \\ \max(\frac{r}{2}; \max_{u \in \mathcal{U}} \rho_0(u) - D(x; u)); & \text{otherwise} \end{cases}$$

and the reward distributions are defined by

$$\Pr(y|x) = p_0(y|x) = \begin{cases} \rho_0(x)r^{1/\epsilon}; & y = r^{-1/\epsilon} \\ 1 - \rho_0(x)r^{1/\epsilon}; & y = 0 \end{cases}$$

The advantage of this construction of  $I_0$  is that the extent to which any other problem instance  $I_i \in \mathcal{I}$  deviates from  $I_0$  can be controlled as follows. Let  $S_i = B(u_i)$

## E. Proof of Lemma 1

**Lemma 1** With probability at least  $1 - \frac{\delta}{T}$ , for all rounds  $t \in [T]$  and all active arms  $x \in A_t$ , we have

$$|b_t(x) - (x)| \leq r_{t+1}(x):$$

*Proof.* Fix  $t \in [T]$  and  $x \in A_t$ . By Lemma 1 in Bubeck et al. (2013), with probability at least  $1 - \frac{\delta}{T}$  the following holds

$$|b_t(x) - (x)| \leq 4^{-\frac{1}{1+\epsilon}} \frac{\log(T^2)}{n_t(x)^{\frac{\epsilon}{1+\epsilon}}} \leq 4^{-\frac{1}{1+\epsilon}} \frac{\log(T^2)}{n_t(x)^{\frac{\epsilon}{1+\epsilon}}} = r_{t+1}(x):$$

Taking the union bound over  $x \in A_t$  and  $t = 1; 2; \dots; T$  and noticing  $|A_t| \leq T$ , we conclude the proof.

## F. Proof of Lemma 2

**Lemma 2** With probability at least  $1 - \frac{\delta}{T}$ , for all rounds  $t \in [T]$  and all active arms  $x \in A_t$ , we have

$$\Delta(x) \leq 3^{-\frac{\epsilon}{2}} r_{t+1}(x):$$

*Proof.* Fix  $t \in [T]$ . For each active arm  $x \in A_t$ , there exist three different scenarios as follows.

(i)  $x$  is pulled by Step 7 of Algorithm 2 in round  $t$ . In this scenario, on one hand, we have  $n_t(x) = 1$  and

$$r_{t+1}(x) = 4^{-\frac{1}{1+\epsilon}} \frac{\log(T^2)}{n_t(x)^{\frac{\epsilon}{1+\epsilon}}} \leq 4^{-\frac{1}{1+\epsilon}} \frac{1}{3^{\frac{\epsilon}{2}}}$$

where we use the fact that  $\log(T^2) \leq \log 4 + 1$  for  $T > 1$ . On the other hand, let  $x_* \in \arg \max_{x \in \mathcal{X}} (x)$  be an optimal arm. We have

$$\Delta(x) = (x_*) - (x) \leq D(x_*, x) \leq 1$$

where the first inequality holds since  $(\cdot)$  is Lipschitz, and the second inequality is due to the assumption in (2). Thus, we obtain  $3^{-\frac{\epsilon}{2}} r_{t+1}(x) \leq \Delta(x)$ .

(ii)  $x$  is pulled by Step 9 of Algorithm 2 in round  $t$ . In this case, we have  $t \geq 2$  and  $n_{t-1}(x) \geq 1$  and the arm selection rule implies

$$b_{t-1}(x) + 2r_t(x) \leq b_{t-1}(x') + 2r_t(x'); \quad \forall x' \in A_t: \quad (25)$$

Note that  $A_t = A_{t-1}$ . By Lemma 1, we get

$$(x) \leq b_{t-1}(x) + r_t(x); \quad b_{t-1}(x') \leq (x') + r_t(x'); \quad \forall x' \in A_t: \quad (26)$$

Combining (25) and (26), we obtain

$$(x) + 3r_t(x) \leq (x') + r_t(x'); \quad \forall x' \in A_t: \quad (27)$$

Note that the execution of Step 9 implies  $X \in [x \in A_t B(x; r_t(x))]$ . Thus, for the optimal arm  $x_*$  there must exist an active arm  $\bar{x}_* \in A_t$  such that

$$D(x_*, \bar{x}_*) \leq r_t(\bar{x}_*)$$

which, together with the Lipschitz property of  $(\cdot)$ , indicates

$$(x_*) \leq (\bar{x}_*) + r_t(\bar{x}_*):$$

Combining the above inequality and (27) with substitution  $x' = \bar{x}_*$ , we obtain

$$(x) + 3r_t(x) \leq (x_*):$$

On the other hand, we have

$$\frac{r_{t+1}(x)}{r_t(x)} = \frac{n_{t-1}(x)}{n_t(x)^{\frac{\epsilon}{1+\epsilon}}} = \frac{n_{t-1}(x)}{n_{t-1}(x) + 1} \leq \frac{1}{2}.$$

Therefore, we get  $3^{-\frac{\epsilon}{2}} r_{t+1}(x) \leq 3r_t(x) \leq \Delta(x)$ :

(iii)  $x$  is not played in round  $t$ . In this scenario, let  $s$  be the last round in which  $x$  is pulled. Then, we have  $r_{t+1}(x) = r_{s+1}(x)$  and the proof reduces to (i) or (ii).

### G. Proof of Lemma 3

**Lemma 3** With probability at least  $1 - 2^{-i}$ , for all  $i \geq 2$ ,

$$|j \bar{A}_T(i)| \leq N_z(2^{-i});$$

*Proof.* By the definition of the  $r$ -zooming number, for  $i = 0, 1, 2, \dots$ , the set  $\bar{A}_T(i)$  can be covered by not more than  $N_z(2^{-i})$  balls of radius at most  $\frac{1}{18 \times 2^i}$ . In the following, we show that each of these balls contains at most one arm from  $\bar{A}_T(i)$ . In fact, suppose that there exist two arms  $u, v \in \bar{A}_T(i)$  falling into the same ball. On one hand, we have

$$D(u; v) \leq \frac{1}{9 \cdot 2^i}; \quad (28)$$

On the other hand, without loss of generality, we assume arm  $u$  is added into the active arm set  $A_T$  before arm  $v$ . Let  $t$  be the time when arm  $v$  is added into  $A_T$ . The execution of Algorithm 2 ensures  $t \geq 2$  and

$$D(u; v) > r_t(u); \quad (29)$$

By Lemma 2, we have

$$r_t(u) = r_{t+1}(u) \geq \frac{\Delta(u)}{3^{t/2}} > \frac{1}{6 \cdot 2^{t/2}} > \frac{1}{9 \cdot 2^i}$$

which, together with (28) and (29), leads to a contradiction. Thus,  $|j \bar{A}_T(i)| \leq N_z(2^{-i})$ .

### H. Proof of Lemma 4

**Lemma 4** With probability at least  $1 - 2^{-i}$ , for all  $i \geq 2$ ,

$$\prod_{x \in \bar{A}_T(i)} n_T(x) \Delta(x) \geq 2^{\frac{i+1}{\epsilon}} \cdot 17^{\frac{\epsilon+1}{\epsilon} - \frac{1}{\epsilon} \log(T^2)} N_z(2^{-i});$$

*Proof.* For any arm  $x \in \bar{A}_T(i)$ , by Lemma 2 we have

$$\Delta(x) \geq 3^{t/2} r_{T+1}(x) \geq 17^{-\frac{1}{1+\epsilon}} \frac{\log(T^2)}{n_T(x)^{\frac{\epsilon}{1+\epsilon}}};$$

Rearranging the above inequality, we obtain

$$n_T(x) \Delta(x) \geq 17^{\frac{\epsilon+1}{\epsilon} - \frac{1}{\epsilon} \log(T^2)} \Delta(x)^{-\frac{1}{\epsilon}} \geq 2^{\frac{i+1}{\epsilon}} \cdot 17^{\frac{\epsilon+1}{\epsilon} - \frac{1}{\epsilon} \log(T^2)}$$

where the second inequality is due to  $\Delta(x) > 2^{-(i+1)}$ : We finish the proof by applying Lemma 3.

### I. Proof of Lemma 11

**Lemma 11** The Kullback–Leibler divergence from  $Q_k$  to  $Q_0$  satisfies

$$KL(Q_0; Q_k) \leq 39 \cdot 200;$$

*Proof.* We first bound the KL divergence from  $p_0$  to  $p_k$ :

$$KL(p_0; p_k) = \int o(x) r^{1/\epsilon} \log \frac{o(x) r^{1/\epsilon}}{k(x) r^{1/\epsilon}} + (1 - \int o(x) r^{1/\epsilon}) \log \frac{1}{1 - \int k(x) r^{1/\epsilon}};$$

In the following, we consider two different scenarios, i.e.,  $x \in X \setminus S_k$  and  $x \in S_k$ .

(i)  $x \in X \setminus S_k$ . By (24), we have  $k(x) = o(x)$  and thus

$$KL(p_0; p_k) = 0; \quad (30)$$

(ii)  $x \geq S_k$ . By (24), we have  $o(x) = k(x) = o(x) + r=8$  which implies

$$\begin{aligned}
 KL(p_0; p_k) &= o(x)r^{1/\epsilon} \log \frac{o(x)r^{1/\epsilon}}{o(x)r^{1/\epsilon}} + (1 - o(x)r^{1/\epsilon}) \log \frac{1 - o(x)r^{1/\epsilon}}{1 - o(x)r^{1/\epsilon}} \\
 &= (1 - o(x)r^{1/\epsilon}) \frac{r^{1+1/\epsilon}=8}{1 - o(x)r^{1/\epsilon} - r^{1+1/\epsilon}=8} \\
 &= (1 - o(x)r^{1/\epsilon} - r^{1+1/\epsilon}=8 + r^{1+1/\epsilon}=8) \frac{r^{1+1/\epsilon}=8}{1 - o(x)r^{1/\epsilon} - r^{1+1/\epsilon}=8} \\
 &= \frac{r^{1+1/\epsilon}}{8} + \frac{r^{1+1/\epsilon}=8^2}{1 - o(x)r^{1/\epsilon} - r^{1+1/\epsilon}=8} \\
 &= \frac{r^{1+1/\epsilon}}{8} + \frac{r^{1+1/\epsilon}=8^2}{1 - 7r^{1+1/\epsilon}=8} \\
 &= \frac{r^{1+1/\epsilon}}{8} + \frac{r^{1+1/\epsilon}=8^2}{4r^{1+1/\epsilon} - 7r^{1+1/\epsilon}=8} \\
 &= \frac{13}{100} r^{1+1/\epsilon}
 \end{aligned} \tag{31}$$

where the second inequality follows from the well-known inequality:  $\log a \leq a - 1$ ;  $8a > 0$ , the third inequality is due to  $o(x) \geq [r=2; 3r=4]$ , and the last inequality holds since  $4r^{1+1/\epsilon} - 4(1=2)^{1+1/\epsilon} - 4(1=2)^2 = 1$ .

We continue the proof of Lemma 11 as follows. Denote by  $KL(\cdot; j)$  the conditional KL divergence also known as conditional relative entropy (Cover & Thomas, 1991; Kleinberg et al., 2013). For  $t = 1; \dots; T$ , we have

$$\begin{aligned}
 KL(Q_0^t; Q_k^t | h^{t-1}) &= \prod_{h^t \in \mathcal{H}^t} Q_0^t(h^t) \log \frac{Q_0^t(h^t | h^{t-1})}{Q_k^t(h^t | h^{t-1})} \\
 &= \prod_{h^t \in \mathcal{H}^t} Q_0^t(h^t) \log \frac{Q_0^t(x_t | h^{t-1})}{Q_k^t(x_t | h^{t-1})} \frac{Q_0^t(y_t | x_t; h^{t-1})}{Q_k^t(y_t | x_t; h^{t-1})} \\
 &= \prod_{h^t \in \mathcal{H}^t} Q_0^t(h^t) \log \frac{Q_0^t(y_t | x_t; h^{t-1})}{Q_k^t(y_t | x_t; h^{t-1})}
 \end{aligned}$$

where the first equality is the definition of conditional KL divergence and the last equality is due to the fact that the distribution of  $x_t$  given  $h^{t-1}$  depends only on the algorithm  $A$ . We proceed as follows

$$\begin{aligned}
 KL(Q_0^t; Q_k^t | h^{t-1}) &= \prod_{h^t \in \mathcal{H}^t} Q_0^t(h^t) \log \frac{Q_0^t(y_t | x_t; h^{t-1})}{Q_k^t(y_t | x_t; h^{t-1})} \\
 &= \prod_{h^{t-1} \in \mathcal{H}^{t-1}} \int_{x_t \in \mathcal{X}} \int_{y_t \in \{0, r-1/\epsilon\}} Q_0^t(y_t | x_t; h^{t-1}) \log \frac{Q_0^t(y_t | x_t; h^{t-1})}{Q_k^t(y_t | x_t; h^{t-1})} d Q_0^t(x_t; h^{t-1}) \\
 &= \prod_{h^{t-1} \in \mathcal{H}^{t-1}} \int_{x_t \in \mathcal{X}} KL(p_0; p_k) d Q_0^t(x_t; h^{t-1}) \\
 &\stackrel{(30)}{=} \prod_{h^{t-1} \in \mathcal{H}^{t-1}} \int_{x_t \in \mathcal{X}} KL(p_0; p_k) d Q_0^t(x_t; h^{t-1}) \\
 &\stackrel{(31)}{=} \prod_{h^{t-1} \in \mathcal{H}^{t-1}} \int_{x_t \in S_k} \frac{13}{100} r^{1+1/\epsilon} d Q_0^t(x_t; h^{t-1}) \\
 &= \frac{13}{100} r^{1+1/\epsilon} Q_0^t(x_t \geq S_k);
 \end{aligned}$$

Finally, by the chain rule of KL divergence we have

$$KL(Q_0; Q_k) = KL(Q_0^T; Q_k^T) = \prod_{t=1}^T KL(Q_0^t; Q_k^t | h^{t-1}) \stackrel{\times}{=} \frac{13}{100} r^{1+1/\epsilon} Q_0^t(x_t \geq S_k) = \frac{13}{100} r^{1+1/\epsilon} \mathbb{E}_{Q_0}[Z_k];$$



where we use the convention that  $h^0 = \emptyset$ . Recalling  $\mathbb{E}_{Q_0}[Z_k] = T/N$  and  $N = \max(2; bTr^{1+1/\epsilon}c)$ , we obtain

$$\mathbb{E}_{Q_0}[Z_k] \leq \frac{3}{2}r^{-(1+1/\epsilon)}$$

which completes the proof.