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

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

the regret of SDTM with parameter r > 0 satisfies

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

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

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

$$(\mathbf{x}_*)$$
 $(\bar{\mathbf{x}}_k)$ $D(\mathbf{x}_*; \bar{\mathbf{x}}_k)$ 2r: (17)

On the other hand, let $\bar{x}_* \ 2 \arg \max_{x_i, i \in [K]} (\bar{x}_i)$ 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 X_* can be upper bounded as follows

$$\mathbb{E} \begin{array}{cccc} & \swarrow & \swarrow & \overset{\#}{} \\ \mathbb{E} & (X_*) & (X_t) & (4 \ T)^{\frac{1}{1+\epsilon}} (16 \ K \log T)^{\frac{\epsilon}{1+\epsilon}} \\ & t=1 & t=1 \end{array}$$
(18)

Combining (17) and (18) and recalling that K

and recalling that $\mathcal{K} = N_c(r)$, we obtain $\mathbb{E} = \begin{pmatrix} X_* \\ t=1 \end{pmatrix} \quad (X_t) = 2rT + (4 T)^{\frac{1}{1+\epsilon}} (16N_c(r)\log T)^{\frac{\epsilon}{1+\epsilon}}$

where we use the fact that (\bar{x}_k) (\bar{X}_*) .

B. Proof of Corollary 2

Corollary 2 We have

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

and thus

$$R(T) \quad O \quad \inf_{\substack{r_0 \in (0,1) \\ \theta \in T \\ \frac{d_z \epsilon + 1}{d_z \epsilon + \epsilon + 1}}} r_0 T + \log T \quad r_0^{-(d_z + 1/\epsilon)}$$

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

Proof. We have

$$\times \underbrace{N_{z}(r)}_{r=2^{-i}:i\in\mathbb{N}, r\geq r_{0}} \underbrace{N_{z}(r)}_{r^{1/\epsilon}} \times \underbrace{2^{id_{z}+i/\epsilon}Z}_{i\in\mathbb{N}:2^{-i}\geq r_{0}} \underbrace{2^{id_{z}+i/\epsilon}Z}_{i=0} \underbrace{2^{id_{z}+i/\epsilon}Z}_{0} \underbrace{2^{id_{z}+i/\epsilon}Z}_{0} \underbrace{2^{id_{z}+i/\epsilon}Z\,\mathrm{d}i}_{0} \underbrace{(2=r_{0})^{d_{z}+1/\epsilon}Z}_{\log(2^{d_{z}+1/\epsilon})} \underbrace{2^{id_{z}+i/\epsilon}Z\,\mathrm{d}i}_{0} \underbrace{(2=r_{0})^{d_{z}+1/\epsilon}Z}_{\log(2^{d_{z}+1/\epsilon})} \underbrace{2^{id_{z}+i/\epsilon}Z\,\mathrm{d}i}_{0} \underbrace{(2=r_{0})^{d_{z}+1/\epsilon}Z}_{\log(2^{d_{z}+1/\epsilon})} \underbrace{2^{id_{z}+i/\epsilon}Z\,\mathrm{d}i}_{0} \underbrace{(2=r_{0})^{d_{z}+1/\epsilon}Z}_{\log(2^{d_{z}+1/\epsilon})} \underbrace{2^{id_{z}+i/\epsilon}Z\,\mathrm{d}i}_{0} \underbrace{2^{id_{z}+i/\epsilon}Z\,\mathrm{d}i}_{0} \underbrace{(2=r_{0})^{d_{z}+1/\epsilon}Z}_{\log(2^{d_{z}+1/\epsilon})} \underbrace{2^{id_{z}+i/\epsilon}Z\,\mathrm{d}i}_{0} \underbrace{2^{id_{z}+i/\epsilon$$

where Z is the zooming constant of (X; D).

C. Proof of Theorem 3

Theorem 3 Assume (2) and (12) hold. With probability at least 1 2, the regret of ADMM satisfies

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

$$R(T)$$
 Θ $T^{\frac{a_z \epsilon + 1}{d_z \epsilon + \epsilon + 1}}$

where d_z is the zooming dimension of $(X; \cdot)$, 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, for all rounds t 2 R and all active arms x 2 A_t , we have

$$ib_t(x)$$
 $(x)j$ $r_{t+1}(x)$:

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

$$jb_t(x) \qquad (x)j \quad (12)^{\frac{1}{1+\epsilon}} \quad \frac{16\log(e^{1/8}T^2 =)}{n_t(x)} \quad \stackrel{\frac{\epsilon}{1+\epsilon}}{=} \quad (12^-)^{\frac{1}{1+\epsilon}} \quad \frac{16\log(e^{1/8}T^2 =)}{n_t(x)} \quad \stackrel{\frac{\epsilon}{1+\epsilon}}{=} r_{t+1}(x):$$

Taking the union bound over $x \ge A_t$ and $t \ge R$ and noticing $jA_t j = T$; $\delta t \ge R$, we conclude the proof.

Lemma 6 With probability at least 1 2, for all rounds $t \ 2[T]$ and all active arms $x \ 2A_t$, we have $\Delta(x) = 3 \frac{D_2}{2} r_{t+1}(x):$

Proof. Fix $t \ge [T]$. For each active arm $x \ge A_t$, 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

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

and hence

$$\begin{split} r_{t+1}(x) &= (12^{-})^{\frac{1}{1+\epsilon}} \quad \frac{16\log\left(e^{1/8}T^{2}=\right)}{n_{t}(x)} \quad \stackrel{\frac{\epsilon}{1+\epsilon}}{(12^{-})^{\frac{1}{1+\epsilon}}} \quad \frac{16\log\left(e^{1/8}T^{2}=\right)}{16\log\left(e^{1/8}T^{2}=\right)+1} \quad \stackrel{\frac{\epsilon}{1+\epsilon}}{(3^{D}\overline{2})^{-\frac{1}{1+\epsilon}}} \quad \frac{35}{36} \quad \frac{\frac{\epsilon}{1+\epsilon}}{36} = \frac{35}{36} \quad 3^{D}\overline{2} \quad \frac{35}{36} \quad -\frac{1}{1+\epsilon}}{36} \quad \frac{35}{36} \quad 3^{D}\overline{2} \quad \frac{35}{36} \quad -\frac{1}{3^{D}\overline{2}}} \end{split}$$

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

(ii) x is pulled by Step 12 of Algorithm 4 in round t. In this case, we have t = 1 2 R and the arm selection rule implies

$$b_{t-1}(x) + 2r_t(x) \quad b_{t-1}(x') + 2r_t(x'); \ 8x' \ 2A_t:$$
 (19)

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

$$(x) \quad b_{t-1}(x) \quad r_t(x); \ b_{t-1}(x') \qquad (x') \quad r_t(x'); \ \partial x' \ 2 \ A_t:$$
(20)

Combining (19) and (20), we obtain

$$(x) + 3r_t(x)$$
 $(x') + r_t(x'); 8x' 2A_t$: (21)

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

$$D(\mathbf{X}_*; \mathbf{\overline{X}}_*) = \mathbf{r}_t(\mathbf{\overline{X}}_*)$$

which, together with the Lipschitz property of , indicates

$$(X_*)$$
 $(\overline{X}_*) + r_t(\overline{X}_*)$:

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

$$(\mathbf{X}) + 3\mathbf{r}_t(\mathbf{X}) \qquad (\mathbf{X}_*)$$

On the other hand, we have

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

Therefore, we get $3^{D_{\overline{2}}} r_{t+1}(x) = 3r_t(x) = \Delta(x)$:

(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, for all $i = 0; 1; 2; \ldots$,

$$|\bar{A}_T(i)| = N_z(2^{-i})$$
:

Lemma 8 With probability at least 1 2, for all $i = 0, 1, 2, \ldots$, $X = n_T(x)\Delta(x) = 2^{\frac{i+1}{\epsilon}} (51^-)^{\frac{1}{\epsilon}} 68\log(e^{1/8}T^2 =)N_z(2^{-i}):$ $x \in \mathcal{A}_T(i)$

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 2(0,1]. Define $=\frac{2^{1/\epsilon} \cdot \epsilon}{\log 2}$ and

$$R_{c}(T) = \inf_{r_{0} \in (0,1)} r_{0}T + \log T \underset{r=2^{-i}: i \in \mathbb{N}, r > r_{0}}{\times} \frac{N_{c}(r)}{r^{1/\epsilon}}$$

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

(i) for each problem instance I 2 I, define

$$R_{z}(T) = \inf_{r_{0} \in (0,1)} r_{0}T + \log T \times \frac{N_{z}(r)}{r^{-2^{-i}:i \in \mathbb{N}, r > r_{0}}} \frac{N_{z}(r)}{r^{1/\epsilon}}$$

in which $N_z(r)$ is the *r*-zooming number of *I*. We have $R_z(T) = 3R = (8 \log T)$: (ii) for any algorithm *A*, there exists at least one problem instance *I* 2 / on which the expected regret of *A* satisfies $\mathbb{E}[R(T)] = R = (2560 \log T)$:

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 /) 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, $R_c(T) = R'_c(T)$.

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

$$R_{c}(T) = \inf_{r_{0} \in (0,1)} r_{0}T + \log T \underset{r=2^{-i}: i \in \mathbb{N}, r \ge r_{0}}{\times} \frac{N_{c}(r)}{r^{1/\epsilon}} \inf_{r_{0} \in (0,1)} r_{0}T + \log T N_{c}(r_{0}) \underset{r=2^{-i}: i \in \mathbb{N}, r \ge r_{0}}{\times} r^{-\frac{1}{\epsilon}}$$

in which the last term can be upper bounded as follows

$$\times_{r=2^{-i}:i\in\mathbb{N},r\geq r_{0}} r^{-\frac{1}{\epsilon}} = \frac{\begin{bmatrix} \log_{2}\frac{1}{r_{0}} \end{bmatrix}}{\sum_{i=0}^{2} 2^{\frac{i}{\epsilon}}} \quad \frac{Z \log_{2}\frac{1}{r_{0}}+1}{2} 2^{\frac{i}{\epsilon}} dI = \frac{(2=r_{0})^{1/\epsilon}}{\log(2^{1/\epsilon})} \quad \frac{2^{1/\epsilon}}{\log 2} r_{0}^{1/\epsilon}:$$

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

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

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

$$N_p(r) = \max f j S j$$
: S is an r-packing of X g:

We have

$$r < 1=2$$
 and $N = N_p(r)$:

Proof. Define function $f(r) = \frac{N_c(r)}{r^{1+1/\epsilon}}$. Since f(1) = 1, $\lim_{r \to 0} f(r) = +7$, and f(r) is decreasing on (0,1), there must exist b 2(0,1) such that f(b) = T (b=2). From the first inequality, we obtain

$$R \quad R_c(T) \qquad \mathbf{b}T + \frac{N_c(\mathbf{b})}{\mathbf{b}^{1/\epsilon}} \log T \qquad \mathbf{b} \ T(1 + \log T)$$

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

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

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

The above lemma ensures that we can find a set of arms $U = fu_1$; $u_N g = X$ such that $\inf_{x,y \in U} D(x,y) = r$. Based on U, we construct a set of problem instances $I = fI_1$; $i \in I_N g$. Let us fix $i \geq [N]$ and describe the construction of I_i : the expected reward function i is defined as

and the reward distributions are defined by

$$\Pr(y|x) = \rho_i(y|x) = \begin{cases} i(x)r^{1/\epsilon}; & y = r^{-1/\epsilon} \\ 1 & i(x)r^{1/\epsilon}; & y = 0 \end{cases}$$
(23)

One can show that for $i = 1; 2; ...; N_i$ is Lipschitz and the (1 +)-th moment of p_i is upper bounded by 7=8.

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

$$R_{z}(T) = \inf_{r_{0} \in (0,1)} r_{0}T + \log T \frac{X}{\rho = 2^{-i} : i \in \mathbb{N}, \rho \ge r_{0}} \frac{N_{z}()}{1/\epsilon} - \frac{3}{4}rT + \log T \frac{X}{\rho = 2^{-i} : i \in \mathbb{N}, \rho \ge \frac{3}{4}r} \frac{N_{z}()}{1/\epsilon} - \frac{3R}{8 \log T}$$

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 $_0$ is defined as

$$_{0}(x) = \begin{pmatrix} \frac{3r}{4}; & x = u_{j}; j \ 2 [N] \\ \max(\frac{r}{2}; \max_{u \in \mathcal{U}} \ _{0}(u) \quad D(x; u)); & \text{otherwise} \end{cases}$$

and the reward distributions are defined by

$$\Pr(y|x) = p_0(y|x) = \begin{pmatrix} 0 & (x)r^{1/\epsilon} \\ 0 & (x)r^{1/\epsilon} \end{pmatrix}, \quad y = r^{-1/\epsilon} \\ 1 & 0 & (x)r^{1/\epsilon} \end{pmatrix}, \quad y = 0$$

,

The advantage of this construction of I_0 is that the extent to which any other problem instance $I_i 2$ / 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 2, for all rounds $t \ge [T]$ and all active arms $x \ge A_t$, we have $jb_t(x) = (x)j = r_{t+1}(x)$:

Proof. Fix $t \ge [T]$ and $x \ge A_t$. By Lemma 1 in Bubeck et al. (2013), with probability at least $1 = \frac{2\delta}{T^2}$ the following holds

$$jb_t(x) \qquad (x)j \quad 4 \quad \frac{1}{1+\epsilon} \quad \frac{\log\left(T^2\right)}{n_t(x)} \quad 4^{-\frac{1}{1+\epsilon}} \quad \frac{\log\left(T^2\right)}{n_t(x)} \quad 4^{-\frac{1}{1+\epsilon}} \quad \frac{\log\left(T^2\right)}{n_t(x)} \quad \epsilon = r_{t+1}(x)$$

Taking the union bound over $x \ge A_t$ and $t = 1; 2; \ldots; T$ and noticing $jA_t j = T; \ 8t \ge [T]$, we conclude the proof.

F. Proof of Lemma 2

Lemma 2 With probability at least 1 2, for all rounds $t \ge [T]$ and all active arms $x \ge A_t$, we have $\Delta(x) = 3 \frac{D}{2} r_{t+1}(x)$:

Proof. Fix $t \ge [T]$. For each active arm $x \ge A_t$, 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

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

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

$$\Delta(\mathbf{x}) = (\mathbf{x}_*) \quad (\mathbf{x}) \quad D(\mathbf{x}_*; \mathbf{x}) \quad 1$$

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

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

$$b_{t-1}(x) + 2r_t(x) \quad b_{t-1}(x') + 2r_t(x'); \ 8x' \ 2A_t:$$
 (25)

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

$$(x) \quad b_{t-1}(x) \quad r_t(x); \ b_{t-1}(x') \qquad (x') \quad r_t(x'); \ 8x' \ 2A_t:$$
(26)

Combining (25) and (26), we obtain

$$(x) + 3r_t(x)$$
 $(x') + r_t(x'); 8x' 2A_t:$ (27)

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

 $D(\mathbf{x}_*; \bar{\mathbf{x}}_*) = \mathbf{r}_t(\bar{\mathbf{x}}_*)$

which, together with the Lipschitz property of , indicates

$$(X_*)$$
 $(\overline{X}_*) + r_t(\overline{X}_*)$:

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

$$(\mathbf{X}) + 3\mathbf{r}_t(\mathbf{X}) \qquad (\mathbf{X}_*)$$

On the other hand, we have

$$\frac{r_{t+1}(x)}{r_t(x)} = \frac{n_{t-1}(x)}{n_t(x)} \stackrel{\frac{\epsilon}{1+\epsilon}}{=} \frac{n_{t-1}(x)}{n_{t-1}(x)+1} \stackrel{\frac{\epsilon}{1+\epsilon}}{\to} \frac{1}{2}:$$

Therefore, we get $3^{P_{\overline{2}}}r_{t+1}(x) = 3r_t(x) = \Delta(x)$:

(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, for all $i \ge \mathbb{N}$,

 $j\bar{A}_{T}(i)j = N_{z}(2^{-i}):$

Proof. By the definition of the *r*-zooming number, for i = 0/1/2/..., 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 $u/2 \bar{A}_T(i)$ falling into the same ball. On one hand, we have

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

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 = 2 and

$$D(u;v) > r_t(u):$$
⁽²⁹⁾

By Lemma 2, we have

$$r_t(u) = r_{t+1}(u) - \frac{\Delta(u)}{3^{p}\overline{2}} > \frac{1}{6^{p}\overline{2}-2^i} > \frac{1}{9-2^i}$$

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

H. Proof of Lemma 4

Lemma 4 With probability at least 1 = 2, for all $i \ge \mathbb{N}$,

$$\times \underset{x \in \mathcal{A}_T(i)}{\times} n_T(x) \Delta(x) \quad 2^{\frac{i+1}{\epsilon}} \quad 17^{\frac{\epsilon+1}{\epsilon} - \frac{1}{\epsilon}} \log \left(T^2 = \right) N_z(2^{-i}):$$

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

$$\Delta(x) = 3^{\mathcal{D}} \overline{2} r_{T+1}(x) = 17^{-\frac{1}{1+\epsilon}} = \frac{\log(T^{2} =)}{n_T(x)} \stackrel{\epsilon}{\longrightarrow}$$

Rearranging the above inequality, we obtain

$$n_T(x)\Delta(x) = 17^{\frac{\epsilon+1}{\epsilon} - \frac{1}{\epsilon}} \log \left(T^2 = \right) \Delta(x)^{-\frac{1}{\epsilon}} = 2^{\frac{\epsilon+1}{\epsilon}} = 17^{\frac{\epsilon+1}{\epsilon} - \frac{1}{\epsilon}} \log \left(T^2 = \right)$$

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

$$KL(Q_0; Q_k) = 39=200$$
:

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

$$\mathcal{K}L(p_0; p_k) = {}_0(x)r^{1/\epsilon}\log - \frac{{}_0(x)r^{1/\epsilon}}{{}_k(x)r^{1/\epsilon}} + (1 - {}_0(x)r^{1/\epsilon})\log - \frac{1 - {}_0(x)r^{1/\epsilon}}{1 - {}_k(x)r^{1/\epsilon}} :$$

In the following, we consider two different scenarios, i.e., $x \ge X$ S_k and $x \ge S_k$. (i) $x \ge X$ S_k . By (24), we have $_k(x) = _0(x)$ and thus

$$\mathcal{K}L(p_0; p_k) = 0. \tag{30}$$

(ii) $x \ge S_k$. By (24), we have $_0(x) = _k(x) = _0(x) + r = 8$ which implies

$$\begin{aligned}
\mathcal{K}L(p_{0};p_{k}) & _{0}(x)r^{1/\epsilon}\log - \frac{_{0}(x)r^{1/\epsilon}}{_{0}(x)r^{1/\epsilon}} + (1 \quad _{0}(x)r^{1/\epsilon})\log - \frac{1 \quad _{0}(x)r^{1/\epsilon}}{1 \quad _{0}(x)r^{1/\epsilon}}r^{1+1/\epsilon}=8 \\
& (1 \quad _{0}(x)r^{1/\epsilon}) - \frac{r^{1+1/\epsilon}=8}{1 \quad _{0}(x)r^{1/\epsilon}}r^{1+1/\epsilon}=8 \\
& = (1 \quad _{0}(x)r^{1/\epsilon} - r^{1+1/\epsilon}=8 + r^{1+1/\epsilon}=8) - \frac{r^{1+1/\epsilon}=8}{1 \quad _{0}(x)r^{1/\epsilon}}r^{1+1/\epsilon}=8 \\
& = \frac{r^{1+1/\epsilon}}{8} + \frac{r^{1+1/\epsilon}=8}{1 \quad _{0}(x)r^{1/\epsilon}}r^{1+1/\epsilon}=8 \\
& = \frac{r^{1+1/\epsilon}}{8} + \frac{r^{1+1/\epsilon}=8}{1 \quad _{0}(x)r^{1/\epsilon}}r^{1+1/\epsilon}=8 \\
& = \frac{r^{1+1/\epsilon}}{8} + \frac{r^{1+1/\epsilon}=8}{2}r^{1+1/\epsilon}=8 \\
& = \frac{r^{1+1/\epsilon}}{8} + \frac{r^{1+1/\epsilon}=8}{2}r^{1+1/\epsilon}=8 \\
& = \frac{13}{100}r^{1+1/\epsilon}
\end{aligned}$$
(31)

where the second inequality follows from the well-known inequality: $\log a = 1$; 8a > 0, the third inequality is due to $_0(x) \ge [r=2; 3r=4]$, 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(; j) the conditional KL divergence also known as conditional relative entropy (Cover & Thomas, 1991; Kleinberg et al., 2013). For t = 1; z : z : T, we have

$$\begin{aligned} \mathcal{K}L(\mathcal{Q}_{0}^{t};\mathcal{Q}_{k}^{t}j\,h^{t-1}) &= \overset{\times}{\underset{h^{t}\in \ ^{t}}{\overset{t}{\longrightarrow}}} \mathcal{Q}_{0}^{t}(h^{t})\log \quad \frac{\mathcal{Q}_{0}^{t}(h^{t}\,j\,h^{t-1})}{\mathcal{Q}_{k}^{t}(h^{t}\,j\,h^{t-1})} \\ &= \overset{\times}{\underset{h^{t}\in \ ^{t}}{\overset{t}{\longrightarrow}}} \mathcal{Q}_{0}^{t}(h^{t})\log \quad \frac{\mathcal{Q}_{0}^{t}(x_{t}\,j\,h^{t-1})}{\mathcal{Q}_{k}^{t}(x_{t}\,j\,h^{t-1})} \quad \frac{\mathcal{Q}_{0}^{t}(y_{t}\,j\,x_{t};h^{t-1})}{\mathcal{Q}_{k}^{t}(y_{t}\,j\,x_{t};h^{t-1})} \\ &= \overset{\times}{\underset{h^{t}\in \ ^{t}}{\overset{t}{\longrightarrow}}} \mathcal{Q}_{0}^{t}(h^{t})\log \quad \frac{\mathcal{Q}_{0}^{t}(y_{t}\,j\,x_{t};h^{t-1})}{\mathcal{Q}_{k}^{t}(y_{t}\,j\,x_{t};h^{t-1})} \end{aligned}$$

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

$$\begin{split} \mathcal{K}L(\mathcal{Q}_{0}^{t};\mathcal{Q}_{k}^{t}\,j\,h^{t-1}) &= \begin{array}{c} \stackrel{\times}{\times} & \mathcal{Q}_{0}^{t}(h^{t})\log & \frac{\mathcal{Q}_{0}^{t}(y_{t}\,j\,x_{t}\,;h^{t-1})}{\mathcal{Q}_{k}^{t}(y_{t}\,j\,x_{t}\,;h^{t-1})} \\ &= \begin{array}{c} & \mathcal{I} & \stackrel{\times}{\times} & \mathcal{I} \\ & \stackrel{h^{t-1}\in & t-1}{\times} & \mathcal{I} \\ & \stackrel{h^{t-1}\in & t-1}{\times} & \mathcal{I} \\ & \stackrel{\star}{\times} \\ & \stackrel{\star}{\times} & \mathcal{I} \\ & \stackrel{\star}{\times} & \mathcal{I} \\ & \stackrel{\star}{\times} & \mathcal{I} \\ & \stackrel{\star}{\times} \\$$

Finally, by the chain rule of KL divergence we have

$$\mathcal{K}L(Q_0;Q_k) = \mathcal{K}L(Q_0^T;Q_k^T) = \underbrace{\overset{\mathcal{K}}{\underset{t=1}{\overset{t}{1}{\overset{t=1}{\overset{t=1}{\overset{t}1}{\overset{t}{1}{\overset{t}{1}{\overset{t}{$$

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

which completes the proof.