

# A PROOFS

## A.1 Proof of Theorem 3.1

*Proof.* To prove the theorem, we need to show that Algorithm 1 will take the same action as the index policy for each state  $s \in \mathcal{S}$  that:

$$\arg \max_{a_i \in \mathcal{A}} q_t(s, a_i) = \arg \max_{a_i \in \mathcal{A}} \hat{z}_t(s^i, a_i).$$

By subtracting a constant from the left-hand side, we get the following equivalent statement:

$$\arg \max_{a_i \in \mathcal{A}} \left( q_t(s, a_i) - \hat{v}_{t+1}(s) \right) = \arg \max_{a_i \in \mathcal{A}} \hat{z}_t(s^i, a_i).$$

Informally, the term  $q_t(s, a_i) - \hat{v}_{t+1}(s)$  measures the advantage of pulling an arm  $a_i$  in comparison with a virtual action of “pulling no arm” instead. Next we show that  $q_t(s, a_i) - \hat{v}_{t+1}(s) = \hat{z}_t(s^i, a_i)$ .

Fix an action  $a_i$ , let  $s_\tau$  be some state, and let  $S_{\tau+1}$  to be the random variable representing the state that follows  $s_\tau^i$  after taking action  $a_i$ . Then:

$$\begin{aligned} q_t(s_\tau, a_i) - \hat{v}_{t+1}(s_\tau) &= \mathbb{E} [\hat{v}_{t+1}(S_{\tau+1}) + r(s_\tau, a_i, S_{\tau+1})] - \hat{v}_{t+1}(s_\tau) = \\ &= \mathbb{E} [r(s_\tau, a_i, S_{\tau+1})] + \mathbb{E} [\hat{v}_{t+1}(S_{\tau+1})] - \hat{v}_{t+1}(s_\tau) = \\ &= r(s_\tau^i, a_i) + \mathbb{E} [\hat{v}_{t+1}(S_{\tau+1})] - \hat{v}_{t+1}(s_\tau) \end{aligned}$$

Using the assumed linearly separable form of  $\hat{v}_t(s_\tau) = \sum_{a_i \in \mathcal{A}} \hat{v}_t^i(s_\tau^i)$ , we get:

$$\begin{aligned} r(s_\tau^i, a_i) + \mathbb{E} [\hat{v}_{t+1}(S_{\tau+1})] - \hat{v}_{t+1}(s_\tau) &= r(s_\tau^i, a_i) + \mathbb{E} \left[ \sum_{a_j \in \mathcal{A}} \hat{v}_{t+1}^j(S_{\tau+1}^j) \right] - \sum_{a_j \in \mathcal{A}} \hat{v}_{t+1}^j(s_\tau^j) = \\ &= r(s_\tau^i, a_i) + \sum_{a_j \in \mathcal{A}} \left( \mathbb{E} [\hat{v}_{t+1}^j(S_{\tau+1}^j)] - \hat{v}_{t+1}^j(s_\tau^j) \right) \end{aligned}$$

Notice next that  $S_{\tau+1}^j = s_\tau^j$  whenever  $a_j \neq a_i$ , since the state of an arm does not change unless the arm is pulled. Thus we can further simplify the sum as follows:

$$\begin{aligned} &r(s_\tau^i, a_i) + \sum_{a_j \in \mathcal{A}} \left( \mathbb{E} [\hat{v}_{t+1}^j(S_{\tau+1}^j)] - \hat{v}_{t+1}^j(s_\tau^j) \right) = \\ &= r(s_\tau^i, a_i) + \sum_{a_j \in \mathcal{A} \setminus \{a_i\}} \left( \mathbb{E} [\hat{v}_{t+1}^j(S_{\tau+1}^j)] - \hat{v}_{t+1}^j(s_\tau^j) \right) + \mathbb{E} [\hat{v}_{t+1}^i(S_{\tau+1}^i)] - \hat{v}_{t+1}^i(s_\tau^i) = \\ &= r(s_\tau^i, a_i) + \sum_{a_j \in \mathcal{A} \setminus \{a_i\}} \left( \mathbb{E} [\hat{v}_{t+1}^j(s_\tau^j)] - \hat{v}_{t+1}^j(s_\tau^j) \right) + \mathbb{E} [\hat{v}_{t+1}^i(S_{\tau+1}^i)] - \hat{v}_{t+1}^i(s_\tau^i) = \\ &= r(s_\tau^i, a_i) + \mathbb{E} [\hat{v}_{t+1}^i(S_{\tau+1}^i)] - \hat{v}_{t+1}^i(s_\tau^i). \end{aligned}$$

Finally, substituting (5) into the equation above gives us:

$$q_t(s_\tau, a_i) - \hat{v}_{t+1}(s_\tau) = \hat{z}_t(s_\tau, a_i),$$

which proves the theorem. □

## A.2 Proof of Lemma 4.1

*Proof.* First, we will use the following decomposition of the regret for a given parameter  $\mu$ :

$$\begin{aligned}
\text{Regret}(\pi, T, \mu) &= \sum_{t=1}^T \mathbb{E} [\mu_{i_u^*} - \mu_{\pi(S_t)}] = \\
&\stackrel{(a)}{=} \sum_{t=1}^T (\mathbb{E} [\mu_{i_u^*} - (q_t(S_t, \pi(S_t)) - v_{t+1}(S_t))] + \mathbb{E} [(q_t(S_t, \pi(S_t)) - v_{t+1}(S_t)) - \mu_{\pi(S_t)}]) = \\
&\stackrel{(b)}{\leq} \sum_{t=1}^T (\mathbb{E} [\mu_{i_u^*} - (q_t(S_t, a_{i_u^*}) - v_{t+1}(S_t))] + \mathbb{E} [(q_t(S_t, \pi(S_t)) - v_{t+1}(S_t)) - \mu_{\pi(S_t)}]) = \\
&= \sum_{t=1}^T \mathbb{E} \left[ \underbrace{(\mu_{i_u^*} - q_t(S_t, a_{i_u^*}) + v_{t+1}(S_t))}_{-\varphi_t(\mu, S_t, a_{i_u^*})} \right] + \sum_{t=1}^T \mathbb{E} \left[ \underbrace{(q_t(S_t, \pi(S_t)) - v_{t+1}(S_t) - \mu_{\pi(S_t)})}_{\varphi_t(\mu, S_t, \pi(S_t))} \right] = \\
&= \sum_{t=1}^T \mathbb{E}_{S_t} [\varphi_t(\mu, S_t, \pi(S_t))] - \sum_{t=1}^T \mathbb{E}_{S_t} [\varphi_t(\mu, S_t, a_{i_u^*})]
\end{aligned}$$

The equality (a) follows by simply adding  $0 = v_{t+1}(S_t) - v_{t+1}(S_t)$ , and the inequality (b) follows from the optimality of  $\pi(S_t)$  with respect to  $v_t$ :

$$q_t(S_t, \pi(S_t)) \geq q_t(S_t, a),$$

for any action  $a \in \mathcal{S}$ . □

## A.3 Proof of Theorem 4.1

The theorem follows directly from the following two lemmas.

**Lemma A.1.** For  $\alpha > 1$  and  $\varphi_t$  computed for  $v^{UCB}$  the following lower bound holds true:

$$\sum_{t=1}^T \mathbb{E}_{S_t} [\varphi_t(\mu, S_t, a_{i_u^*})] \geq - \sum_{t=1}^T \frac{1}{t^{2\alpha-1}} \geq O(1).$$

*Proof.*

$$\begin{aligned}
\mathbb{E}_{S_t} [\varphi_t(\mu, S_t, a_{i_u^*})] &= \mathbb{E}_{S_t} [q_t(S_t, a_{i_u^*}) - v_{t+1}(S_t) - \mu_{i_u^*}] = \mathbb{E} [z^{UCB}(S_t, a_{i_u^*}) - \mu_{i_u^*}] = \\
&= \mathbb{E} \left[ r(S_t, a_{i_u^*}) + \sqrt{\frac{\alpha \log t}{T_{i_u^*}^{i_u^*}(S_t^{i_u^*}, a_{i_u^*})}} - \mu_{i_u^*} \right] \geq \\
&\stackrel{(a)}{\geq} -\mathbb{P} \left[ r(S_t, a_{i_u^*}) - \mu_{i_u^*} \leq \sqrt{\frac{\alpha \log t}{T_{i_u^*}^{i_u^*}(S_t^{i_u^*}, a_{i_u^*})}} \right] \geq \\
&\stackrel{(b)}{\geq} -t \exp \left( -2T_{i_u^*}^{i_u^*}(S_t^{i_u^*}, a_{i_u^*}) \frac{\alpha \log t}{T_{i_u^*}^{i_u^*}(S_t^{i_u^*}, a_{i_u^*})} \right) = \\
&= -t \exp(-2\alpha \log t) = -\frac{1}{t^{2\alpha-1}}
\end{aligned}$$

where (a) follows by rewriting expectation and the fact that the rewards are bounded between 0 and 1, and (b) follows from Hoeffding's inequality and the union bound over possible values of  $T_{i_u^*}$ . Summing the value above for  $t = 1 \dots T$  proves the lemma. □

**Lemma A.2.** Assume  $\alpha > 1$  and fix an action  $a_i \in \mathcal{A}$ . Then for a policy  $\pi$  greedy with respect to  $v^{UCB}$  for following inequality holds true:

$$\sum_{t=1}^T \mathbb{E}_{S_t} [\varphi_t(\mu, S_t, \pi(S_t))] \leq O(\sqrt{T \log T})$$

*Proof.* Consider the bound for a single action:

$$\begin{aligned} \mathbb{E}_{S_t} [\varphi_t(\mu, S_t, a_i)] &= \mathbb{E}_{S_t} [q_t(S_t, a_i) - v_{t+1}(S_t) - \mu_i] = \mathbb{E} [z^{UCB}(S_t, a_i) - \mu_{i_t}^*] \\ &= \mathbb{E} \left[ r(S_t, a_i) + \sqrt{\frac{\alpha \log t}{T_i(S_t^i, a_i)}} - \mu_i \right] = \\ &= \mathbb{E} \left[ r(S_t, a_i) - \sqrt{\frac{\alpha \log t}{T_i(S_t^i, a_i)}} + 2\sqrt{\frac{\alpha \log t}{T_i(S_t^i, a_i)}} - \mu_i \right] = \\ &= \mathbb{E} \left[ r(S_t, a_i) - \sqrt{\frac{\alpha \log t}{T_i(S_t^i, a_i)}} - \mu_i \right] + \mathbb{E} \left[ 2\sqrt{\frac{\alpha \log t}{T_i(S_t^i, a_i)}} \right] = \\ &\stackrel{(a)}{\leq} \mathbb{P} \left[ r(S_t, a_i) - \mu_i \leq \sqrt{\frac{\alpha \log t}{T_i(S_t^i, a_i)}} \right] + \mathbb{E} \left[ 2\sqrt{\frac{\alpha \log t}{T_i(S_t^i, a_i)}} \right] \\ &\stackrel{(b)}{\leq} \frac{1}{t^{2\alpha-1}} + \mathbb{E} \left[ 2\sqrt{\frac{\alpha \log t}{T_i(S_t^i, a_i)}} \right] \end{aligned}$$

The equalities above follow readily using algebra; (a) and (b) follow by an argument that is identical to the proof of Lemma A.1.

Considering just the state in which arm  $a_i$  is pulled, it remains to upper bound the following expression:

$$\begin{aligned} \sum_{t=1}^T 2 \mathbb{E} \left[ \mathbf{1}_{a_i=\pi(S_t)} \sqrt{\frac{\alpha \log t}{T_i(S_t^i, a_i)}} \right] &\leq \sum_{t=1}^T 2 \mathbb{E} \left[ \mathbf{1}_{a_i=\pi(S_t)} \sqrt{\frac{\alpha \log T}{T_i(S_t^i, a_i)}} \right] \leq \\ &\stackrel{(d)}{\leq} \sum_{t=1}^T 2 \mathbb{E} \left[ \sqrt{\frac{\alpha \log T}{t}} \right] \stackrel{(e)}{\leq} O(\sqrt{T \log T}). \end{aligned}$$

The inequality (d) follows by upper bounding the error by assuming that the arm was pulled in every step, and (e) follows by upper bounding the sum by an integral. See, for example, the proof of Proposition 2 in Russo and Van Roy [2014].  $\square$