
Supplemental Materials for *Dynamic Interruption Policies for Reinforcement Learning Game Playing Using Multi-Sampling Multi-Armed Bandits*

A. Appendix

A.1. Proofs of the following two equations introduced in Section 2.2.

$$\begin{aligned}
 \mathbb{E}_{\lambda_{[1:\infty]}} [C_{I_t,t}] &= \mathbb{E}_{\lambda} \left[\sum_J \left(\prod_{j=1}^{J-1} \epsilon_j \right) \cdot c_{I_t,t,J} \right] \\
 &= \mathbb{E}_c \left[\sum_J \left(\prod_{j=1}^{J-1} (1 - c_{I_t,t,j}) \right) c_{I_t,t,J} \right] \\
 &= 1,
 \end{aligned} \tag{1}$$

$$\begin{aligned}
 \mathbb{E}_{\lambda, I_t} \left[\frac{R_{I_t,t} \mathbb{1}_{I_t=i}}{p_{i,t}} \right] &= \sum_{I_t=1}^K p_{I_t,t} \mathbb{E} \left[\frac{R_{I_t,t} \mathbb{1}_{I_t=i}}{p_{i,t}} \right] = \mathbb{E}[R_{i,t}] \\
 &= \mathbb{E}_{\lambda} \left[\sum_J \left(\prod_{j=1}^{J-1} \epsilon_j \right) \cdot (r_{i,t,J}) \right] \\
 &= \mathbb{E}_{r,c} \left[\sum_J \left(\prod_{j=1}^{J-1} (1 - c_{i,t,j}) \right) r_{i,t,J} \right] \\
 &= \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]},
 \end{aligned} \tag{2}$$

A.2. Proof of Theorem 1.

We first proof a lemma described as follows,

Lemma 1: Let

$$\beta = \sqrt{\frac{\ln(K)}{nK}}, \tilde{R}_{i,t} = \frac{R_{I_t,t} \mathbb{1}_{I_t=i} + 3\beta}{p_{i,t}}, n \geq 100,$$

we can obtain the following inequality,

$$\mathbb{E}_{\lambda, I_t} \left[\max_i \left(\beta \sum_{t=1}^n \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \sum_{t=1}^n \tilde{R}_{i,t} \right) \right] \leq \ln(K) + 3. \tag{3}$$

Proof of Lemma 1: Let $\beta = 0.33\sqrt{\frac{\ln(K)}{nK}}$, $B = 2 \ln(n+1)$. It is easy to verify that,

$$\mathbb{1}_{C_{I_t,t} > B} \leq \frac{\exp(\frac{C_{I_t,t}}{2}) - 1}{\exp(\frac{B}{2}) - 1} = \frac{\exp(\frac{C_{I_t,t}}{2}) - 1}{n}. \tag{4}$$

Using (4) and the inequality $\exp(x) \leq 1 + x + x^2$ for $x \leq 1$, we have,

$$\begin{aligned}
 \mathbb{P}(C_{I_t,t} > B) &= \sum_J \left(\prod_{j=1}^{J-1} (1 - c_{I_t,t,j}) \right) c_{I_t,t,J} \cdot \mathbb{1}_{C_{I_t,t} > B} \\
 &\leq \frac{1}{n} \sum_J \left(\prod_{j=1}^{J-1} (1 - c_{I_t,t,j}) \right) c_{I_t,t,J} \cdot \left(\exp \left(\sum_{j=1}^J \frac{c_{I_t,t,j}}{2} \right) - 1 \right) \\
 &= \frac{1}{n} \left(\frac{\mathbb{E}[c_{I_t,t} \exp(\frac{c_{I_t,t}}{2})]}{1 - \mathbb{E}[(1 - c_{I_t,t}) \exp(\frac{c_{I_t,t}}{2})]} - 1 \right) \\
 &\leq \frac{1}{n} \left(\frac{\mathbb{E}[c_{I_t,t} + \frac{c_{I_t,t}^2}{2} + \frac{c_{I_t,t}^3}{4}]}{\mathbb{E}[\frac{c_{I_t,t}}{2} + \frac{c_{I_t,t}^2}{4} + \frac{c_{I_t,t}^3}{4}]} - 1 \right) \\
 &< \frac{1}{n},
 \end{aligned}$$

where the last inequality comes from the fact $\frac{\mathbb{E}[c_{I_t,t} + \frac{c_{I_t,t}^2}{2} + \frac{c_{I_t,t}^3}{4}]}{\mathbb{E}[\frac{c_{I_t,t}}{2} + \frac{c_{I_t,t}^2}{4} + \frac{c_{I_t,t}^3}{4}]} < 2$.

Moreover,

$$\begin{aligned}
 \mathbb{E}_{\lambda} \left[\frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} \right]^2 &= \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} \mathbb{E}_c \left[\sum_J \left(\prod_{j=1}^{J-1} (1 - c_{I_t,t,j}) \right) c_{I_t,t,J} \cdot \left(\sum_{j=1}^J c_{I_t,t,j} \right)^2 \right] \\
 &= \left(\frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} \right)^2 \times \left(2 - \frac{\mathbb{E}[c_{I_t,t}^2]}{\mathbb{E}[c_{I_t,t}]} \right).
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 \mathbb{E}_{\lambda, I_t} \left[\frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} \cdot \frac{R_{I_t,t} \mathbb{1}_{I_t=i}}{p_{i,t}} \right] &= \left(\frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} \right)^2 \times \left(2 - \frac{\mathbb{E}[c_{i,t}^2]}{\mathbb{E}[c_{i,t}]} \right), \\
 \mathbb{E}_{\lambda, I_t} \left[\frac{R_{I_t,t} \mathbb{1}_{I_t=i}}{p_{i,t}} \right]^2 &= \frac{1}{p_{i,t}} \left(2 \left(\frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} \right)^2 + \frac{\mathbb{E}[r_{i,t}]^2}{\mathbb{E}[c_{i,t}]} - 2 \left(\frac{\mathbb{E}[r_{i,t}] \mathbb{E}[r_{i,t} c_{i,t}]}{\mathbb{E}[c_{i,t}]^2} \right) \right).
 \end{aligned}$$

And thus,

$$\mathbb{E}_{\lambda, I_t} \left[\beta \frac{E[r_{i,t}]}{E[c_{i,t}]} C_{I_t,t} - \beta \frac{R_{I_t,t} \mathbb{1}_{I_t=i}}{p_{i,t}} \right]^2 \leq \frac{3\beta^2}{p_{i,t}}.$$

When $C_{I_t,t} \leq B$, we have $\beta \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} \leq 1$. Then,

$$\begin{aligned}
 & \mathbb{E}_{\lambda, I_t} \exp \left(\beta \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \tilde{R}_{I_t,t} \right) \\
 &= \mathbb{P}(C_{I_t,t} \leq B) \cdot \mathbb{E}_{\lambda, I_t} \left[\exp \left(\beta \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \tilde{R}_{I_t,t} \right) \middle| C_{I_t,t} \leq B \right] \\
 &+ \mathbb{P}(C_{I_t,t} > B) \cdot \mathbb{E}_{\lambda, I_t} \left[\exp \left(\beta \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \tilde{R}_{I_t,t} \right) \middle| C_{I_t,t} > B \right] \\
 &\leq \mathbb{P}(C_{I_t,t} \leq B) \cdot \mathbb{E}_{\lambda, I_t} \left[1 + \mathbb{E}_{\lambda, I_t} \left[\beta \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \frac{R_{I_t,t} \mathbb{1}_{I_t=i}}{p_{i,t}} \right] + \mathbb{E}_{\lambda, I_t} \left[\beta \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \frac{R_{I_t,t} \mathbb{1}_{I_t=i}}{p_{i,t}} \right]^2 \middle| C_{I_t,t} \leq B \right] \\
 &\times \exp \left(-\frac{3\beta^2}{p_{i,t}} \right) + \mathbb{P}(C_{I_t,t} > B) \cdot \left[\exp \left(\beta \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} (B+1) \right) \middle| C_{I_t,t} > B \right] \\
 &\leq \left(1 + \mathbb{E}_{\lambda, I_t} \left[\beta \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \frac{R_{I_t,t} \mathbb{1}_{I_t=i}}{p_{i,t}} \right] + \mathbb{E}_{\lambda, I_t} \left[\beta \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \frac{R_{I_t,t} \mathbb{1}_{I_t=i}}{p_{i,t}} \right]^2 \right) \\
 &\times \exp \left(-\frac{3\beta^2}{p_{i,t}} \right) + \frac{\exp(\beta(B+1))}{n} \\
 &\leq \left(1 + \frac{3\beta^2}{p_{i,t}} \right) \times \exp \left(-\frac{3\beta^2}{p_{i,t}} \right) + \frac{2}{n} \\
 &\leq 1 + \frac{2}{n}.
 \end{aligned}$$

Thus, we have,

$$\begin{aligned}
 & \mathbb{E}_{\lambda, I_t} \left[\exp \left(\beta \sum_{t=1}^n \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \sum_{t=1}^n \tilde{R}_{I_t,t} \right) \right] \\
 &\leq \left(1 + \frac{2}{n} \right)^n \\
 &\leq e^2.
 \end{aligned}$$

Moreover, Markov's inequality implies $\mathbb{P}(X > \ln(\delta^{-1})) \leq \delta \mathbb{E}e^X$ and thus, with probability at least $1 - \delta e^2$,

$$\beta \sum_{t=1}^n \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \sum_{t=1}^n \tilde{R}_{I_t,t} \leq \ln(\delta^{-1}).$$

As a consequence, with probability at most $K\delta e^2$,

$$\max_i \left(\beta \sum_{t=1}^n \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \sum_{t=1}^n \tilde{R}_{I_t,t} \right) > \ln(\delta^{-1}).$$

This equals to say, with probability at most δ ,

$$\max_i \left(\beta \sum_{t=1}^n \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \sum_{t=1}^n \tilde{R}_{I_t,t} \right) - \ln K - 2 > \ln(\delta^{-1}).$$

Using following equation:

$$\mathbb{E}[W] \leq \int_0^1 \frac{1}{\delta} \mathbb{P}(W > \ln(\frac{1}{\delta})) d\delta.$$

In particular, taking

$$W = \max_i \left(\beta \sum_{t=1}^n \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \beta \sum_{t=1}^n \tilde{R}_{I_t,t} \right) - \ln K - 2,$$

yields $\mathbb{E}[W] \leq 1$, which is equivalent to inequality (3).

Proof of Theorem 1: One can immediately see that

$$\begin{aligned}
 \mathbb{E}[\text{Reg}(n)] &= \mathbb{E} \max_{\lambda, I_t} \left[\gamma \sum_{t=1}^n \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} + (1-\gamma) \left(\sum_{t=1}^n \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \sum_{t=1}^n \tilde{R}_{i,t} + \sum_{t=1}^n \tilde{R}_{i,t} \right) - \sum_{t=1}^n R_{I_t,t} \right] \\
 &\leq \gamma \mathbb{E} \sum_{\lambda, I_t} C_{I_t,t} + (1-\gamma) \max_i \left(\sum_{t=1}^n \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} C_{I_t,t} - \sum_{t=1}^n \tilde{R}_{i,t} \right) + \mathbb{E} \max_{\lambda, I_t} \left[(1-\gamma) \sum_{t=1}^n \tilde{R}_{i,t} - \sum_{t=1}^n R_{I_t,t} \right] \\
 &\leq \gamma n + \frac{(1-\gamma)(3+\ln K)}{\beta} + \mathbb{E} \max_{\lambda, I_t} \left[(1-\gamma) \sum_{t=1}^n \tilde{R}_{i,t} - \sum_{t=1}^n R_{I_t,t} \right].
 \end{aligned}$$

Where the last inequality comes from Lemma 1.

Let $u = (\frac{1}{K}, \dots, \frac{1}{K})$ be the uniform distribution over the arms, $w_t = \frac{p_t - u\gamma}{1-\gamma}$ be the distribution induced by Exp3.PMS at time t without the mixing. Since,

$$w_{i,1} = 0, \quad w_{i,t} = \frac{\exp(\eta \sum_{s=1}^{t-1} \tilde{R}_{i,s})}{\sum_{k=1}^K \exp(\eta \sum_{s=1}^{t-1} \tilde{R}_{k,s})},$$

we have:

$$\begin{aligned}
 \sum_{t=1}^n \ln \mathbb{E}_{i \sim w_t} \exp(\eta \tilde{R}_{i,t}) &= \sum_{t=1}^n \ln \sum_{i=1}^K \frac{\exp(\eta \sum_{\tau=1}^{t-1} \tilde{R}_{i,\tau}) \exp(\eta \tilde{R}_{i,t})}{\sum_{k=1}^K \exp(\eta \sum_{\tau=1}^{t-1} \tilde{R}_{k,\tau})} \\
 &= \ln \prod_{t=1}^n \frac{\sum_{k=1}^K \exp(\eta \sum_{\tau=1}^t \tilde{R}_{k,\tau})}{\sum_{k=1}^K \exp(\eta \sum_{\tau=1}^{t-1} \tilde{R}_{k,\tau})} \\
 &= \ln \left(\sum_{k=1}^K \exp(\eta \sum_{t=1}^n \tilde{R}_{k,t}) \right) - \ln(K) \\
 &\geq \max_k \ln \left(\exp(\eta \sum_{t=1}^n \tilde{R}_{k,t}) \right) - \ln(K) \\
 &= \max_k \eta \sum_{t=1}^n \tilde{R}_{k,t} - \ln(K).
 \end{aligned} \tag{5}$$

Additionally, it is easily to verify that,

$$- \sum_{t=1}^n R_{I_t,t} = - \sum_{t=1}^n \mathbb{E}_{i \sim p_t} \tilde{R}_{i,t} + 3\beta n K.$$

Then,

$$\begin{aligned}
 \mathbb{E}[\text{Reg}(n)] &\leq \gamma n + \frac{(1-\gamma)(3+\ln K)}{\beta} + \mathbb{E} \max_{\lambda, I_t} \left[(1-\gamma) \sum_{t=1}^n \tilde{R}_{i,t} - \sum_{t=1}^n R_{I_t,t} \right] \\
 &\leq \gamma n + \frac{(1-\gamma)(3+\ln K)}{\beta} + \mathbb{E}_{\lambda, I_t} \left[\frac{1-\gamma}{\eta} \sum_{t=1}^n \ln \mathbb{E}_{i \sim w_t} \exp(\eta \tilde{R}_{i,t}) \right. \\
 &\quad \left. - \sum_{t=1}^n \mathbb{E}_{i \sim p_t} \tilde{R}_{i,t} + 3\beta n K \right] + \frac{(1-\gamma) \ln K}{\eta} \\
 &= \mathbb{E}_{\lambda, I_t} \left[(1-\gamma) \left(\frac{1}{\eta} \ln \mathbb{E}_{i \sim w_t} \exp(\eta \tilde{R}_{i,t}) - \mathbb{E}_{k \sim p_t} \eta \tilde{R}_{k,t} - \gamma \mathbb{E}_{i \sim \mu} \tilde{R}_{i,t} \right) \right. \\
 &\quad \left. + \gamma n + \frac{(1-\gamma)(3+\ln K)}{\beta} + 3\beta n K + \frac{(1-\gamma) \ln K}{\eta} \right] \\
 &< \mathbb{E}_{\lambda, I_t} \left[\frac{1-\gamma}{\eta} (\ln \mathbb{E}_{i \sim w_t} \exp(\eta \tilde{R}_{i,t}) - \mathbb{E}_{k \sim p_t} \eta \tilde{R}_{k,t}) \right] \\
 &\quad + \gamma n + \frac{(1-\gamma)(3+\ln K)}{\beta} + 3\beta n K + \frac{(1-\gamma) \ln K}{\eta}.
 \end{aligned}$$

Using the inequalities $\ln x \leq x - 1$ and $\exp(x) \leq 1 + x + x^2$, for all $x \leq 1$, as well as the fact that $(1 + 3\beta)\eta K \leq \gamma$:

$$\begin{aligned}
 \mathbb{E}_\lambda \left[\ln \mathbb{E}_{i \sim w_t} \exp(\eta \tilde{R}_{i,t}) - \mathbb{E}_{k \sim p_t} \eta \tilde{R}_{k,t} \right] &\leq \mathbb{E}_\lambda \mathbb{E}_{i \sim w_t} [\exp(\eta \tilde{R}_{i,t}) - 1 - \eta \tilde{R}_{i,t}] \\
 &= \mathbb{E}_\lambda \mathbb{E}_{i \sim w_t} \left[\exp\left(\frac{\eta R_{i,t} \mathbb{1}_{I_t=i} + 3\eta\beta}{p_{i,t}}\right) - 1 - \eta \tilde{R}_{i,t} \right] \\
 &= w_{I_t,t} \frac{\mathbb{E}[c_{I_t,t} \exp(\frac{\eta r_{I_t,t} + 3\eta\beta}{p_{I_t,t}})]}{1 - \mathbb{E}[(1 - c_{I_t,t}) \exp(\frac{\eta r_{I_t,t}}{p_{I_t,t}})]} + \sum_{i \neq I_t} w_{i,t} \exp\left(\frac{3\eta\beta}{p_{i,t}}\right) \\
 &\quad - 1 - \frac{w_{I_t,t} \eta \mathbb{E}[r_{I_t,t}]}{p_{I_t,t} \mathbb{E}[c_{I_t,t}]} - \mathbb{E}_{i \sim w_t} \frac{3\eta\beta}{p_{i,t}} \\
 &= w_{I_t,t} \left(\frac{\mathbb{E}[c_{I_t,t} \exp(\frac{\eta r_{I_t,t} + 3\eta\beta}{p_{I_t,t}})]}{1 - \mathbb{E}[(1 - c_{I_t,t}) \exp(\frac{\eta r_{I_t,t}}{p_{I_t,t}})]} - \frac{\eta \mathbb{E}[r_{I_t,t}]}{p_{I_t,t} \mathbb{E}[c_{I_t,t}]} \right. \\
 &\quad \left. - \frac{3\eta\beta}{p_{I_t,t}} - 1 \right) + \sum_{i \neq I_t} w_{i,t} \left(\exp\left(\frac{3\eta\beta}{p_{i,t}}\right) - \frac{3\eta\beta}{p_{i,t}} - 1 \right) \\
 &\triangleq w_{I_t,t} A_1 + \sum_{i \neq I_t} w_{i,t} A_2.
 \end{aligned} \tag{6}$$

For convenience, we omit the index I_t and t without ambiguity, and denote $x = \frac{\eta r_{I_t,t}}{p_{I_t,t}}$, $b = \frac{3\eta\beta}{p_{I_t,t}}$, $c = c_{I_t,t}$. We rewrite the term A_1 to:

$$\begin{aligned}
 A_1 &= \frac{\mathbb{E}[ce^{x+b}]}{1 - \mathbb{E}[(1-c)e^x]} - \frac{\mathbb{E}[x]}{\mathbb{E}[c]} - b - 1 \\
 &= \frac{\mathbb{E}[c]\mathbb{E}[ce^{x+b}] + \mathbb{E}[(1-c)e^x](\mathbb{E}[x] + b\mathbb{E}[c] + \mathbb{E}[c])}{\mathbb{E}[c](1 - \mathbb{E}[(1-c)e^x])} - \frac{(\mathbb{E}[x] + b\mathbb{E}[c] + \mathbb{E}[c])}{\mathbb{E}[c](1 - \mathbb{E}[(1-c)e^x])} \\
 &\leq \frac{\mathbb{E}[c]\mathbb{E}[c((x+b)^2 + (x+b) + 1)] - (\mathbb{E}[x] + b\mathbb{E}[c] + \mathbb{E}[c])}{\mathbb{E}[c](1 - \mathbb{E}[(1-c)(x^2 + x + 1)])} + \frac{\mathbb{E}[(1-c)(x^2 + x + 1)](\mathbb{E}[x] + b\mathbb{E}[c] + \mathbb{E}[c])}{\mathbb{E}[c](1 - \mathbb{E}[(1-c)(x^2 + x + 1)])} \\
 &= \frac{(\mathbb{E}[x] + b\mathbb{E}[c])\mathbb{E}[x^2 - x^2c + x + bc] + b\mathbb{E}[xc]\mathbb{E}[c]}{\mathbb{E}[c](1 - \mathbb{E}[(1-c)(x^2 + x + 1)])} - \frac{b\mathbb{E}[x]\mathbb{E}[c] + \mathbb{E}[c]\mathbb{E}[x^2] - \mathbb{E}[x]\mathbb{E}[xc]}{\mathbb{E}[c](1 - \mathbb{E}[(1-c)(x^2 + x + 1)])} \\
 &\leq \frac{(\mathbb{E}[x] + b\mathbb{E}[c])\mathbb{E}[x^2 - x^2c + x + bc] + \mathbb{E}[c]\mathbb{E}[x^2]}{\mathbb{E}[c](1 - \mathbb{E}[(1-c)(x^2 + x + 1)])} \\
 &\leq \frac{(\mathbb{E}[x] + b\mathbb{E}[c])\mathbb{E}[x^2 - x^2c + x + bc + \mathbb{E}[c]\frac{\eta}{p_{I_t,t}}]}{\mathbb{E}[c](1 - \mathbb{E}[(1-c)(x^2 + x + 1)])} \\
 &\leq \frac{(\mathbb{E}[x] + b\mathbb{E}[c])\mathbb{E}[(1-c)cx\frac{\eta}{p_{I_t,t}} + x + bc + \mathbb{E}[c]\frac{\eta}{p_{I_t,t}}]}{\mathbb{E}[c](1 - \mathbb{E}[(1-c)(x^2 + x + 1)])}.
 \end{aligned}$$

Since $(1-c)c \leq 0.25$, we have,

$$\begin{aligned}
 A_1 &\leq \frac{\eta^2}{p_{I_t,t}} \frac{(2.25 + 3\beta)\mathbb{E}[c_{I_t,t}]}{1 - \mathbb{E}[(1 - c_{I_t,t})(1 + \frac{\eta}{p_{I_t,t}}r + (\frac{\eta}{p_{I_t,t}}r)^2)]} \times \left(\frac{\mathbb{E}[r_{I_t,t}]}{\mathbb{E}[c_{I_t,t}]p_{I_t,t}} + \frac{3\beta}{p_{I_t,t}} \right) \\
 &\leq \frac{\eta^2}{p_{I_t,t}} \frac{(2.25 + 3\beta)\mathbb{E}[c_{I_t,t}]}{(1 - \frac{\eta}{p_{I_t,t}})\mathbb{E}[c_{I_t,t}]} \left(\frac{\mathbb{E}[r_{I_t,t}]}{\mathbb{E}[c_{I_t,t}]p_{I_t,t}} + \frac{3\beta}{p_{I_t,t}} \right) \\
 &= \frac{\eta^2}{p_{I_t,t}} \frac{(2.25 + 3\beta)p_{I_t,t}}{p_{I_t,t} - \eta} \left(\frac{\mathbb{E}[r_{I_t,t}]}{\mathbb{E}[c_{I_t,t}]p_{I_t,t}} + \frac{3\beta}{p_{I_t,t}} \right) \\
 &= \frac{\eta^2}{p_{I_t,t}} \frac{(2.25 + 3\beta)\gamma/K}{\gamma/K - \eta} \left(\frac{\mathbb{E}[r_{I_t,t}]}{\mathbb{E}[c_{I_t,t}]p_{I_t,t}} + \frac{3\beta}{p_{I_t,t}} \right).
 \end{aligned} \tag{7}$$

At the same time,

$$\begin{aligned}
 A_2 &= \exp\left(\frac{3\eta\beta}{p_{i,t}}\right) - \frac{3\eta\beta}{p_{i,t}} - 1 \\
 &\leq \left(\frac{3\eta\beta}{p_{i,t}}\right)^2 + \frac{3\eta\beta}{p_{i,t}} + 1 - \frac{3\eta\beta}{p_{i,t}} - 1 \\
 &\leq \frac{\eta^2}{p_{i,t}} \frac{(2.25 + \beta)\gamma/K}{\gamma/K - \eta} \left(\frac{3\beta}{p_{i,t}}\right).
 \end{aligned} \tag{8}$$

Combine (7) and (8), we have

$$\begin{aligned}
 \mathbb{E}_\lambda \left[\ln \mathbb{E}_{i \sim w_t} \exp(\eta \tilde{R}_{i,t}) - \mathbb{E}_{i \sim w_t} \eta \tilde{R}_{i,t} \right] &\leq \sum_{i=1}^K \frac{w_i}{p_i} \eta^2 \frac{(2.25 + 3\beta)\gamma/K}{\gamma/K - \eta} \left(\frac{\mathbb{1}_{i=I_t} + 3\beta}{p_{i,t}} \right) \\
 &\leq \frac{\eta^2}{1-\gamma} \frac{(2.25 + 3\beta)\gamma/K}{\gamma/K - \eta} \sum_{i=1}^K \frac{\mathbb{1}_{i=I_t} + 3\beta}{p_{i,t}} \\
 &= \frac{\eta^2}{1-\gamma} \frac{(2.25 + 3\beta)\gamma/K}{\gamma/K - \eta} \sum_{I_t=1}^K p_{I_t} \sum_{i=1}^K \frac{\mathbb{1}_{i=I_t} + 3\beta}{p_{i,t}} \\
 &= \frac{\eta^2}{1-\gamma} \frac{(2.25 + 3\beta)\gamma/K}{\gamma/K - \eta} \left(K + \sum_{I_t=1}^K \sum_{i=1}^K \frac{3p_{I_t}\beta}{p_{i,t}} \right) \\
 &= \frac{\eta^2}{1-\gamma} \frac{(2.25 + 3\beta)\gamma/K}{\gamma/K - \eta} \left(K + \sum_{i=1}^K \frac{\beta}{p_{i,t}} \right) \\
 &\leq \frac{\eta^2}{1-\gamma} \frac{(2.25 + 3\beta)\gamma/K}{\gamma/K - \eta} \left(K + \sum_{i=1}^K \frac{K\beta}{\gamma} \right) \\
 &\leq \frac{\eta^2}{1-\gamma} \frac{(2.25 + 3\beta)\gamma/K}{\gamma/K - \eta} \left(K + K \frac{K\beta}{\gamma} \right).
 \end{aligned}$$

To sum up,

$$\begin{aligned}
 \mathbb{E}[\text{Reg}(n)] &\leq \gamma n + \frac{(1-\gamma)(3 + \ln K)}{\beta} + \frac{(1-\gamma) \ln K}{\eta} + 3\beta n K + \sum_{t=1}^n \mathbb{E}_\lambda \left[\frac{1-\gamma}{\eta} \left(\ln \mathbb{E}_{i \sim w_t} \exp(\eta \tilde{R}_{i,t}) - \mathbb{E}_{i \sim w_t} \eta \tilde{R}_{i,t} \right) \right] \\
 &\leq \gamma n + \frac{(1-\gamma)(3 + \ln K)}{\beta} + \frac{(1-\gamma) \ln K}{\eta} + 3\beta n K + n \frac{1-\gamma}{\eta} \frac{\eta^2}{1-\gamma} \frac{(2.25 + 3\beta)\gamma/K}{\gamma/K - \eta} \left(K + K \frac{3K\beta}{\gamma} \right) \\
 &\leq \gamma n + \frac{(1-\gamma)(3 + \ln K)}{\beta} + \frac{(1-\gamma) \ln K}{\eta} + 3\beta n K + \eta n K \frac{(2.25 + 3\beta)\gamma/K}{\gamma/K - \eta} \left(1 + \frac{3K\beta}{\gamma} \right) \\
 &\leq 17.33 \sqrt{nK \ln K}.
 \end{aligned}$$

A.3. The details of the modification of the Exp3.PMS algorithm

Algorithm 1 Modification of Exp3.PMS

Parameters: $\eta \in \mathbb{R}^+$ and $\gamma, \beta \in [0, 1]$.

Let p_1 be the uniform distribution over $1, \dots, K$.

For each round $t = 1, \dots, n$

(1) Draw an arm I_t from the probability distribution p_t , and let $C_{I_t,t} = 0, R_{i,t} = 0, a_{i,t} = 0$.

while $C_{I_t,t} < 2 \ln(n+1)$ **do**

 Pull arm I_t and record the reward $r_{I_t,t}$ and the cost $c_{I_t,t}$ ($c_{I_t,t} \leq \tau_{I_t}$),

$$R_{I_t,t} \leftarrow R_{I_t,t} + r_{I_t,t}, \quad C_{I_t,t} \leftarrow C_{I_t,t} + c_{I_t,t}.$$

For each arm $i = 1, \dots, (I_t - 1)$,

$$\begin{aligned} R_{i,t} &\leftarrow R_{i,t} + a_{i,t} \cdot r_{I_t,t} \mathbb{1}_{c_{I_t,t} \leq \tau_i}, \\ a_{i,t} &\leftarrow a_{i,t} \max \left(0, \frac{(\tau_i - c_{I_t,t}) \tau_{I_t}}{(\tau_{I_t} - c_{I_t,t}) \tau_i} \right). \end{aligned}$$

For each arm $i = (I_t + 1), \dots, K$,

$$R_{i,t} = 0.$$

 Break the loop with probability $\frac{c_{I_t,t}}{\tau_{I_t}}$.

end while

(2) Compute the estimated gain for each arm:

$$\tilde{R}_{i,t} = \frac{1}{\tau_i} \left(R_{1,t} + \sum_{j=2}^{\min(I_t, i)} \frac{(R_{j,t} - R_{j-1,t})}{\sum_{k=i}^K p_{k,t}} \right) + \frac{3\beta}{\sum_{k=i}^K p_{k,t}},$$

and update the estimated cumulative gain for each arm:

$$q_{i,t} = \sum_{s=1}^t \tilde{R}_{i,s}. \quad (9)$$

(3) Compute the new probability distribution over the arms $p_{t+1} = (p_{1,t+1}, \dots, p_{K,t+1})$ where:

$$p_{i,t+1} = (1 - \gamma) \frac{\exp(\eta q_{i,t})}{\sum_{i=1}^K \exp(\eta q_{i,t})} + \gamma \mathbb{1}_{i=K}$$

As at each round, some information of arms $1, \dots, I_t - 1$ can also be monitored, we assign a value instead of 0 for these arms. By this way, the variance of the estimate $\tilde{R}_{i,t}$ can be reduced. For the multi-sampling process at round t , after pulling an arm with a cost of $c_{I_t,t}$, the bandit agent will continue to play the current arm with a probability of $1 - \frac{c_{I_t,t}}{\tau_{I_t}}$. This means that if the last game level is not completed but interrupted, then the current arm will be abandoned, and a new round will start. We use $\lambda_j = (r_j, c_j, \epsilon_j)$ denoting a variable vector, where the index j is used to count the number of sampling, and we omit the notations of I_t and t without ambiguity. ϵ_j is a random variable depending on c_j , satisfying that $\mathbb{P}(\epsilon_j = 1) = 1 - \frac{c_j}{\tau_{I_t}}$, and $\mathbb{P}(\epsilon_j = 0) = \frac{c_j}{\tau_{I_t}}$. If we don't take the constraint $C_{I_t,t} < 2 \ln(n+1)$ into account and let the sampling process naturally end, then we can obtain that,

$$\mathbb{E}_{\lambda_{[1:\infty]}} [C_{I_t,t}] = \tau_{I_t}. \quad (10)$$

Let

$$\hat{R}_{i,t} = \frac{1}{\tau_i} \left(R_{1,t} + \sum_{j=2}^{\min(I_t, i)} \frac{(R_{j,t} - R_{j-1,t})}{\sum_{k=i}^K p_{k,t}} \right).$$

Then we can verify that $\hat{R}_{i,t}$ is the unbiased estimator of $\frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]}$.

$$\mathbb{E}_{\lambda, I_t} [\hat{R}_{i,t}] = \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]}.$$
 (11)

Proof: At round t , we pull arm I_t and receive $c_{I_t,t}$ and $r_{I_t,t}$. If instead, we pull arm i ($i < I_t$) and receive $c_{i,t}$ and $r_{i,t}$, then 1) if $c_{I_t,t} > \tau_i$, we have $c_{i,t} = \tau_i$, $r_{i,t} = 0$; 2) if $c_{I_t,t} \leq \tau_i$, we have $c_{i,t} = c_{I_t,t}$, $r_{i,t} = r_{I_t,t}$.

Thus,

$$\begin{aligned} & \mathbb{E}[r_{I_t,t} | c_{I_t,t} \leq \tau_i] \cdot \mathbb{P}(c_{I_t,t} \leq \tau_i) \\ &= \mathbb{E}[r_{i,t} | c_{I_t,t} \leq \tau_i] \cdot \mathbb{P}(c_{I_t,t} > \tau_i) + \mathbb{E}[r_{I_t,t} | c_{I_t,t} \leq \tau_i] \cdot \mathbb{P}(c_{I_t,t} > \tau_i) \\ &= \mathbb{E}[r_i], \end{aligned}$$
 (12)

$$\begin{aligned} & \mathbb{E}\left[\frac{\tau_i - c_{I_t,t}}{\tau_i} \middle| c_{I_t,t} \leq \tau_i\right] \mathbb{P}(c_{I_t,t} \leq \tau_i) \\ &= \mathbb{E}\left[\frac{\tau_i - c_{i,t}}{\tau_i} \middle| c_{I_t,t} \leq \tau_i\right] \mathbb{P}(c_{I_t,t} \leq \tau_i) \\ &= \mathbb{E}\left[\frac{\tau_i - c_{i,t}}{\tau_i} \middle| c_{I_t,t} \leq \tau_i\right] \mathbb{P}(c_{I_t,t} \leq \tau_i) \\ &+ \mathbb{E}\left[\frac{\tau_i - c_{i,t}}{\tau_i} \middle| c_{I_t,t} > \tau_i\right] \mathbb{P}(c_{I_t,t} > \tau_i) \\ &= \mathbb{E}\left[\frac{\tau_i - c_{i,t}}{\tau_i}\right]. \end{aligned}$$
 (13)

Based on (12) and (13),

$$\begin{aligned} \mathbb{E}[R_{i,t}] &= \mathbb{E}(r_{I_t,t,1} \leq \tau_i | c_{I_t,t,1} \leq \tau_i) \mathbb{P}(c_{I_t,t,1} \leq \tau_i) + \sum_{J=2}^{\infty} \left(\mathbb{E}[r_{I_t,t,J} \leq \tau_i | c_{I_t,t,J} \leq \tau_i] \mathbb{P}(c_{I_t,t,J} \leq \tau_i) \right) \\ & \prod_{j=1}^{J-1} \mathbb{E}\left[\frac{\tau_{I_t} - c_{I_t,t,j}}{\tau_{I_t}} \frac{(\tau_i - c_{I_t,t,j})\tau_{I_t}}{(\tau_{I_t} - c_{I_t,t,j})\tau_i} \middle| c_{I_t,t,j} \leq \tau_i\right] \mathbb{P}(c_{I_t,t,j} \leq \tau_i) \\ &= \mathbb{E}[r_{i,t,1}] + \sum_{J=2}^{\infty} \left(\mathbb{E}[r_{i,t,J}] \prod_{j=1}^{J-1} \mathbb{E}\left[\frac{\tau_i - c_{I_t,t,j}}{\tau_i} \middle| c_{I_t,t,j} \leq \tau_i\right] \mathbb{P}(c_{I_t,t,j} \leq \tau_i) \right) \\ &= \mathbb{E}[r_{i,t}] + \mathbb{E}[r_{i,t}] \sum_J \mathbb{E}\left[1 - \frac{c_{i,t}}{\tau_i}\right]^J \\ &= \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} \tau_i. \end{aligned}$$
 (14)

And finally,

$$\begin{aligned} \mathbb{E}_{\lambda, I_t} [\hat{\mu}_{i,t}] &= \frac{1}{\tau_i} \mathbb{E}_{\lambda} \left(\sum_{I_t=1}^K p_{I_t,t} R_{1,t} + \sum_{I_t=1}^K p_{I_t,t} \sum_{j=2}^{\min(i, I_t)} \frac{(R_{j,t} - R_{j-1,t})}{\sum_{k=i}^K p_{k,t}} \right) \\ &= \frac{1}{\tau_i} \left(\frac{\mathbb{E}[r_{1,t}]}{\mathbb{E}[c_{1,t}]} \tau_1 + \sum_{j=2}^i \sum_{I_t=j}^K p_{I_t,t} \frac{(\frac{\mathbb{E}[r_{j,t}]}{\mathbb{E}[c_{j,t}]} \tau_j - \frac{\mathbb{E}[r_{j-1,t}]}{\mathbb{E}[c_{j-1,t}]} \tau_{j-1})}{\sum_{k=i}^K p_{k,t}} \right) \\ &= \frac{\mathbb{E}[r_{i,t}]}{\mathbb{E}[c_{i,t}]} \tau_i. \end{aligned}$$