

## A Proof of core FEWA guarantees

**Lemma 1.** *On the favorable event  $\xi_t$ , if an arm  $i$  passes through a filter of window  $h$  at round  $t$ , i.e.,  $i \in \mathcal{K}_h$ , then the average of its  $h$  last pulls satisfies*

$$\bar{\mu}_i^h(N_{i,t}^{\pi_F}) \geq \mu_t^+(\pi_F) - 4c(h, \delta_t). \quad (5)$$

*Proof.* Let  $i$  be an arm that passed a filter of window  $h$  at round  $t$ . First, we use the confidence bound for the estimates and we pay the cost of keeping all the arms up to a distance  $2c(h, \delta_t)$  of  $\hat{\mu}_{\max,t}^h$ ,

$$\bar{\mu}_i^h(N_{i,t}) \geq \hat{\mu}_i^h(N_{i,t}) - c(h, \delta_t) \geq \hat{\mu}_{\max,t}^h - 3c(h, \delta_t) \geq \max_{i \in \mathcal{K}_h} \bar{\mu}_i^h(N_{i,t}) - 4c(h, \delta_t), \quad (11)$$

where in the last inequality, we used that that for all  $i \in \mathcal{K}_h$ ,

$$\hat{\mu}_{\max,t}^h \geq \hat{\mu}_i^h(N_{i,t}) \geq \bar{\mu}_i^h(N_{i,t}) - c(h, \delta_t).$$

Second, since the means of arms are decaying, we know that

$$\mu_t^+(\pi_F) \triangleq \mu_{i_t^*}(N_{i_t^*,t}) \leq \mu_{i_t^*}(N_{i_t^*,t} - 1) = \bar{\mu}_{i_t^*}^1(N_{i_t^*,t}) \leq \max_{i \in \mathcal{K}} \bar{\mu}_i^1(N_{i,t}) = \max_{i \in \mathcal{K}_1} \bar{\mu}_i^1(N_{i,t}). \quad (12)$$

Third, we show that the largest average of the last  $h'$  means of arms in  $\mathcal{K}_{h'}$  is increasing with  $h'$ ,

$$\forall h' \leq N_{i,t} - 1, \max_{i \in \mathcal{K}_{h'+1}} \bar{\mu}_i^{h'+1}(N_{i,t}) \geq \max_{i \in \mathcal{K}_{h'}} \bar{\mu}_i^{h'}(N_{i,t}).$$

To show the above property, we remark that thanks to our selection rule, the arm that has the largest average of means, always passes the filter. Formally, we show that  $\arg \max_{i \in \mathcal{K}_{h'}} \bar{\mu}_i^{h'}(N_{i,t}) \subseteq \mathcal{K}_{h'+1}$ . Let  $i_{\max}^{h'} \in \arg \max_{i \in \mathcal{K}_{h'}} \bar{\mu}_i^{h'}(N_{i,t})$ . Then for such  $i_{\max}^{h'}$ , we have

$$\hat{\mu}_{i_{\max}^{h'}}^{h'}(N_{i_{\max}^{h'},t}) \geq \bar{\mu}_{i_{\max}^{h'}}^{h'}(N_{i_{\max}^{h'},t}) - c(h', \delta_t) \geq \bar{\mu}_{\max,t}^{h'} - c(h', \delta_t) \geq \hat{\mu}_{\max,t}^{h'} - 2c(h', \delta_t),$$

where the first and the third inequality are due to confidence bounds on estimates, while the second one is due to the definition of  $i_{\max}^{h'}$ .

Since the arms are decaying, the average of the last  $h' + 1$  mean values for a given arm is always greater than the average of the last  $h'$  mean values and therefore,

$$\max_{i \in \mathcal{K}_{h'}} \bar{\mu}_i^{h'}(N_{i,t}) = \bar{\mu}_{i_{\max}^{h'}}^{h'}(N_{i_{\max}^{h'},t}) \leq \bar{\mu}_{i_{\max}^{h'}}^{h'+1}(N_{i_{\max}^{h'},t}) \leq \max_{i \in \mathcal{K}_{h'+1}} \bar{\mu}_i^{h'+1}(N_{i,t}), \quad (13)$$

because  $i_{\max}^{h'} \in \mathcal{K}_{h'+1}$ . Gathering Equations 11, 12, and 13 leads to the claim of the lemma,

$$\bar{\mu}_i^h(N_{i,t}) \stackrel{(11)}{\geq} \max_{i \in \mathcal{K}_h} \bar{\mu}_i^h(N_{i,t}) - 4c(h, \delta_t) \stackrel{(13)}{\geq} \max_{i \in \mathcal{K}_1} \bar{\mu}_i^1(N_{i,t}) - 4c(h, \delta_t) \stackrel{(12)}{\geq} \mu_t^+(\pi_F) - 4c(h, \delta_t). \quad \square$$

**Corollary 1.** *Let  $i \in \text{OP}$  be an arm overpulled by FEWA at round  $t$  and  $h_{i,t} \triangleq N_{i,t}^{\pi_F} - N_{i,t}^{\pi^*} \geq 1$  be the difference in the number of pulls w.r.t. the optimal policy  $\pi^*$  at round  $t$ . On the favorable event  $\xi_t$ , we have*

$$\mu_t^+(\pi_F) - \bar{\mu}_i^{h_{i,t}}(N_{i,t}) \leq 4c(h_{i,t}, \delta_t). \quad (8)$$

*Proof.* If  $i$  was pulled at round  $t$ , then by the condition at Line 10 of Algorithm 1, it means that  $i$  passes through all the filters from  $h = 1$  up to  $N_{i,t}$ . In particular, since  $1 \leq h_{i,t} \leq N_{i,t}$ ,  $i$  passed the filter for  $h_{i,t}$ , and thus we can apply Lemma 1 and conclude

$$\bar{\mu}_i^{h_{i,t}}(N_{i,t}) \geq \mu_t^+(\pi_F) - 4c(h_{i,t}, \delta_t). \quad (14) \quad \square$$

## B Proofs of auxiliary results

**Lemma 2.** Let  $h_{i,t}^\pi \triangleq |N_{i,T}^\pi - N_{i,T}^{\pi^*}|$ . For any policy  $\pi$ , the regret at round  $T$  is no bigger than

$$R_T(\pi) \leq \sum_{i \in \text{OP}} \sum_{h=0}^{h_{i,T}^\pi - 1} \left[ \xi_{t_i^\pi(N_{i,T}^{\pi^*} + h)} \right] \left( \mu_T^+(\pi) - \mu_i(N_{i,T}^{\pi^*} + h) \right) + \sum_{t=0}^T \left[ \bar{\xi}_t \right] Lt.$$

We refer to the first sum above as to  $A_\pi$  and to the second one as to  $B$ .

*Proof.* We consider the regret at round  $T$ . From Equation 3, the decomposition of regret in terms of overpulls and underpulls gives

$$R_T(\pi) = \sum_{i \in \text{UP}} \sum_{t'=N_{i,T}^{\pi^*}+1}^{N_{i,T}^{\pi^*}} \mu_i(t') - \sum_{i \in \text{OP}} \sum_{t'=N_{i,T}^{\pi^*}+1}^{N_{i,T}^{\pi^*}} \mu_i(t').$$

In order to separate the analysis for each arm, we upper-bound all the rewards in the first sum by their maximum  $\mu_T^+(\pi) \triangleq \max_{i \in \mathcal{K}} \mu_i(N_{i,T}^{\pi^*})$ . This upper bound is tight for problem-independent bound because one cannot hope that the unexplored reward would decay to reduce its regret in the worst case. We also notice that there are as many terms in the first double sum (number of underpulls) than in the second one (number of overpulls). This number is equal to  $\sum_{\text{OP}} h_{i,T}^\pi$ . Notice that this does *not* mean that for each arm  $i$ , the number of overpulls equals to the number of underpulls, which cannot happen anyway since an arm cannot be simultaneously underpulled and overpulled. Therefore, we keep only the second double sum,

$$R_T(\pi) \leq \sum_{i \in \text{OP}} \sum_{t'=0}^{h_{i,T}^\pi - 1} \left( \mu_T^+(\pi_F) - \mu_i(N_{i,T}^{\pi^*} + t') \right). \quad (15)$$

Then, we need to separate overpulls that are done under  $\xi_t$  and under  $\bar{\xi}_t$ . We introduce  $t_i^\pi(n)$ , the round at which  $\pi$  pulls arm  $i$  for the  $n$ -th time. We now make the round at which each overpull occurs explicit,

$$\begin{aligned} R_T(\pi) &\leq \sum_{i \in \text{OP}} \sum_{t'=0}^{h_{i,T}^\pi - 1} \sum_{t=0}^T \left[ t_i^\pi(N_{i,T}^{\pi^*} + t') = t \right] \left( \mu_T^+(\pi) - \mu_i(N_{i,T}^{\pi^*} + t') \right) \\ &\leq \underbrace{\sum_{i \in \text{OP}} \sum_{t'=0}^{h_{i,T}^\pi - 1} \sum_{t=0}^T \left[ t_i^\pi(N_{i,T}^{\pi^*} + t') = t \wedge \xi_t \right] \left( \mu_T^+(\pi) - \mu_i(N_{i,T}^{\pi^*} + t') \right)}_{A_\pi} \\ &\quad + \underbrace{\sum_{i \in \text{OP}} \sum_{t'=0}^{h_{i,T}^\pi - 1} \sum_{t=0}^T \left[ t_i^\pi(N_{i,T}^{\pi^*} + t') = t \wedge \bar{\xi}_t \right] \left( \mu_T^+(\pi) - \mu_i(N_{i,T}^{\pi^*} + t') \right)}_B. \end{aligned}$$

For the analysis of the pulls done under  $\xi_t$  we do not need to know at which round it was done. Therefore,

$$A_\pi \leq \sum_{i \in \text{OP}} \sum_{t'=0}^{h_{i,T}^\pi - 1} \left[ \xi_{t(N_{i,T}^{\pi^*} + t')} \right] \left( \mu_T^+(\pi) - \mu_i(N_{i,T}^{\pi^*} + t') \right).$$

For FEWA, it is not easy to directly guarantee the low probability of overpulls (the second sum). Thus, we upper-bound the regret of each overpull at round  $t$  under  $\xi_t$  by its maximum value  $Lt$ . While this is done to ease FEWA analysis, this is valid for any policy  $\pi$ . Then, noticing that we can have at most 1 overpull per round  $t$ , i.e.,  $\sum_{i \in \text{OP}} \sum_{t'=0}^{h_{i,T}^\pi - 1} \left[ t_i^\pi(N_{i,T}^{\pi^*} + t') = t \right] \leq 1$ , we get

$$B \leq \sum_{t=0}^T \left[ \bar{\xi}_t \right] Lt \sum_{i \in \text{OP}} \sum_{t'=0}^{h_{i,T}^\pi - 1} \left[ t_i^\pi(N_{i,T}^{\pi^*} + t') = t \right] \leq \sum_{t=0}^T \left[ \bar{\xi}_t \right] Lt.$$

Therefore, we conclude that

$$R_T(\pi) \leq \underbrace{\sum_{i \in \text{OP}} \sum_{t'=0}^{h_{i,T}^\pi - 1} \left[ \xi_{t'}^\pi(N_{i,t}^* + t') \right] \left( \mu_T^+(\pi) - \mu_i(N_{i,T}^{\pi^*} + t') \right)}_{A_\pi} + \underbrace{\sum_{t=0}^T \left[ \bar{\xi}_t \right] L t}_B.$$

□

**Lemma 3.** Let  $h_{i,t} \triangleq h_{i,t}^{\pi_F} = |N_{i,t}^{\pi_F} - N_{i,t}^{\pi^*}|$ . For policy  $\pi_F$  with parameters  $(\alpha, \delta_0)$ ,  $A_{\pi_F}$  defined in Lemma 2 is upper-bounded by

$$\begin{aligned} A_{\pi_F} &\triangleq \sum_{i \in \text{OP}} \sum_{t'=0}^{h_{i,T} - 1} \left[ \xi_{t'}^{\pi_F}(N_{i,t}^* + t') \right] \left( \mu_T^+(\pi_F) - \mu_i(N_{i,T}^{\pi^*} + t') \right) \\ &\leq \sum_{i \in \text{OP}_\xi} \left( 4\sqrt{2\alpha\sigma^2 \log_+(KT\delta_0^{-1/\alpha})} + 4\sqrt{2\alpha\sigma^2 (h_{i,T}^\xi - 1) \log_+(KT\delta_0^{-1/\alpha})} + L \right). \end{aligned}$$

*Proof.* First, we define  $h_{i,T}^\xi \triangleq \max\{h \leq h_{i,T} \mid \xi_{t'}^{\pi_F}(N_{i,t}^* + h)\}$ , the last overpull of arm  $i$  pulled at round  $t_i \triangleq t_i^{\pi_F}(N_{i,t}^* + h_{i,T}^\xi) \leq T$  under  $\xi_t$ . Now, we upper-bound  $A_{\pi_F}$  by including all the overpulls of arm  $i$  until the  $h_{i,T}^\xi$ -th overpull, even the ones under  $\bar{\xi}_t$ ,

$$A_{\pi_F} \triangleq \sum_{i \in \text{OP}} \sum_{t'=0}^{h_{i,T}^{\pi_F} - 1} \left[ \xi_{t'}^{\pi_F}(N_{i,t}^* + t') \right] \left( \mu_T^+(\pi_F) - \mu_i(N_{i,T}^{\pi^*} + t') \right) \leq \sum_{i \in \text{OP}_\xi} \sum_{t'=0}^{h_{i,T}^\xi - 1} \left( \mu_T^+(\pi_F) - \mu_i(N_{i,T}^{\pi^*} + t') \right),$$

where  $\text{OP}_\xi \triangleq \{i \in \text{OP} \mid h_{i,T}^\xi \geq 1\}$ . We can therefore split the second sum of  $h_{i,T}^\xi$  term above into two parts. The first part corresponds to the first  $h_{i,T}^\xi - 1$  (possibly zero) terms (overpulling differences) and the second part to the last  $(h_{i,T}^\xi - 1)$ -th one. Recalling that at round  $t_i$ , arm  $i$  was selected under  $\xi_{t_i}$ , we apply Corollary 1 to bound the regret caused by previous overpulls of  $i$  (possibly none),

$$A_{\pi_F} \leq \sum_{i \in \text{OP}_\xi} \mu_T^+(\pi_F) - \mu_i(N_{i,T}^* + h_{i,T}^\xi - 1) + 4(h_{i,T}^\xi - 1)c(h_{i,T}^\xi - 1, \delta_{t_i}) \quad (16)$$

$$\leq \sum_{i \in \text{OP}_\xi} \mu_T^+(\pi_F) - \mu_i(N_{i,T}^* + h_{i,T}^\xi - 1) + 4(h_{i,T}^\xi - 1)c(h_{i,T}^\xi - 1, \delta_T) \quad (17)$$

$$\leq \sum_{i \in \text{OP}_\xi} \mu_T^+(\pi_F) - \mu_i(N_{i,T}^* + h_{i,T}^\xi - 1) + 4\sqrt{2\alpha\sigma^2 (h_{i,T}^\xi - 1) \log_+(KT\delta_0^{-1/\alpha})}, \quad (18)$$

with  $\log_+(x) \triangleq \max(\log(x), 0)$ . The second inequality is obtained because  $\delta_t$  is decreasing and  $c(\cdot, \cdot, \delta)$  is decreasing as well. The last inequality is the definition of confidence interval in Proposition 4 with  $\log_+(KT^\alpha) \leq \alpha \log_+(KT)$  for  $\alpha > 1$ . If  $N_{i,T}^{\pi^*} = 0$  and  $h_{i,T}^\xi = 1$  then

$$\mu_T^+(\pi_F) - \mu_i(N_{i,T}^{\pi^*} + h_{i,T}^\xi - 1) = \mu^+(\pi_F) - \mu_i(0) \leq L,$$

since and  $\mu^+(\pi_F) \leq L$  and  $\mu_i(0) \geq 0$  by the assumptions of our setting. Otherwise, we can decompose

$$\mu_T^+(\pi_F) - \mu_i(N_{i,T}^{\pi^*} + h_{i,T}^\xi - 1) = \underbrace{\mu_T^+(\pi_F) - \mu_i(N_{i,T}^{\pi^*} + h_{i,T}^\xi - 2)}_{A_1} + \underbrace{\mu_i(N_{i,T}^{\pi^*} + h_{i,T}^\xi - 2) - \mu_i(N_{i,T}^{\pi^*} + h_{i,T}^\xi - 1)}_{A_2}.$$

For term  $A_1$ , since arm  $i$  was overpulled at least once by FEWA, it passed at least the first filter. Since this  $h_{i,T}^\xi$ -th overpull is done under  $\xi_{t_i}$ , by Lemma 1 we have that

$$A_1 \leq 4c(1, \delta_{t_i}) \leq 4c(1, K^{-1}T^{-\alpha}) \leq 4\sqrt{2\alpha\sigma^2 \log_+(KT\delta_0^{-1/\alpha})}.$$

The second difference,  $A_2 = \mu_i(N_{i,T}^{\pi^*} + h_{i,T}^\xi - 2) - \mu_i(N_{i,T}^{\pi^*} + h_{i,T}^\xi - 1)$  cannot exceed  $L$ , since by the assumptions of our setting, the maximum decay in one round is bounded. Therefore, we further upper-bound Equation 18 as

$$A_{\pi_F} \leq \sum_{i \in \text{OP}_\xi} \left( 4\sqrt{2\alpha\sigma^2 \log_+ \left( KT\delta_0^{-1/\alpha} \right)} + 4\sqrt{2\alpha\sigma^2 \left( h_{i,T}^\xi - 1 \right) \log_+ \left( KT\delta_0^{-1/\alpha} \right)} + L \right). \quad (19)$$

□

**Lemma 4.** Let  $\zeta(x) = \sum_n n^{-x}$ . Thus, with  $\delta_t = \delta_0/(Kt^\alpha)$  and  $\alpha > 4$ , we can use Proposition 4 and get

$$\mathbb{E}[B] \triangleq \sum_{t=0}^T p(\xi_t) Lt \leq \sum_{t=0}^T \frac{Lt\delta_0}{2t^{\alpha-2}} \leq L\delta_0 \frac{\zeta(\alpha-3)}{2}.$$

## C Minimax regret analysis of FEWA

**Theorem 1.** For any rotating bandit scenario with means  $\{\mu_i(n)\}_{i,n}$  satisfying Asm. 1 with bounded decay  $L$  and any time horizon  $T$ , FEWA run with  $\alpha = 5$  and  $\delta_t = 1/(Kt^5)$ , suffers an expected regret<sup>7</sup> of

$$\mathbb{E}[R_T(\pi_F)] \leq 13\sigma(\sqrt{KT} + K)\sqrt{\log(KT)} + KL.$$

*Proof.* To get the problem-independent upper bound for FEWA, we need to upper-bound the regret by quantities which do not depend on  $\{\mu_i\}_i$ . The proof is based on Lemma 2, where we bound the expected values of terms  $A_{\pi_F}$  and  $B$  from the statement of the lemma. We start by noting that on high-probability event  $\xi_T$ , we have by Lemma 3 and  $\alpha = 5$  that

$$A_{\pi_F} \leq \sum_{i \in \text{OP}_\xi} \left( 4\sqrt{10\sigma^2 \log(KT)} + 4\sqrt{10\sigma^2(h_i - 1) \log(KT)} + L \right).$$

Since  $\text{OP}_\xi \subseteq \text{OP}$  and there are at most  $K - 1$  overpulled arms, we can upper-bound the number of terms in the above sum by  $K - 1$ . Next, the total number of overpulls  $\sum_{i \in \text{OP}} h_{i,T}$  cannot exceed  $T$ . As square-root function is concave we can use Jensen's inequality. Moreover, we can deduce that the worst allocation of overpulls is the uniform one, i.e.,  $h_{i,T} = T/(K - 1)$ ,

$$\begin{aligned} A_{\pi_F} &\leq (K - 1)(4\sqrt{10\sigma^2 \log(KT)} + L) + 4\sqrt{10\sigma^2 \log(KT)} \sum_{i \in \text{OP}} \sqrt{(h_{i,T} - 1)} \\ &\leq (K - 1)(4\sqrt{10\sigma^2 \log(KT)} + L) + 4\sqrt{10\sigma^2(K - 1)T \log(KT)}. \end{aligned} \quad (20)$$

Now, we consider the expectation of term  $B$  from Lemma 2. According to Lemma 4, with  $\alpha = 5$  and  $\delta_0 = 1$ ,

$$\mathbb{E}[B] \leq \frac{L\zeta(2)}{2} = \frac{L\pi^2}{12}. \quad (21)$$

Therefore, using Lemma 2 together with Equations 20 and 21, we bound the total expected regret as

$$\mathbb{E}[R_T(\pi_F)] \leq 4\sqrt{10\sigma^2(K - 1)T \log(KT)} + (K - 1)(4\sqrt{10\sigma^2 \log(KT)} + L) + \frac{L\pi^2}{6}. \quad (22)$$

□

**Corollary 3.** FEWA run with  $\alpha > 3$  and  $\delta_0 \triangleq 2\delta/\zeta(\alpha - 2)$  achieves with probability  $1 - \delta$ ,

$$R_T(\pi_F) = A_{\pi_F} \leq 4\sqrt{2\alpha\sigma^2 \log_+ \left( \frac{KT}{\delta_0^{1/\alpha}} \right)} \left( K - 1 + \sqrt{(K - 1)T} \right) + (K - 1)L.$$

<sup>7</sup>See Corollary 3 and 4 for the high-probability result.

*Proof.* We consider the event  $\bigcup_{t \leq T} \xi_t$  which happens with probability

$$1 - \sum_{t \leq T} \frac{Kt^2 \delta_t}{2} \leq 1 - \sum_{t \leq T} \frac{Kt^2 \delta_t}{2} \leq 1 - \frac{\zeta(\alpha - 2)\delta_0}{2}.$$

Therefore, by setting  $\delta_0 \triangleq 2\delta/\zeta(\alpha - 2)$ , we have that  $B = 0$  with probability  $1 - \delta$  since  $\lceil \xi_t \rceil = 0$  for all  $t$ . We can then use the same analysis of  $A_{\pi_F}$  as in Theorem 1 to get

$$R_T(\pi_F) = A_{\pi_F} \leq 4 \sqrt{2\alpha\sigma^2 \log_+ \left( \frac{KT}{\delta_0^{1/\alpha}} \right)} \left( K - 1 + \sqrt{(K - 1)T} \right) + (K - 1)L.$$

□

## D Problem-dependent regret analysis of FEWA

**Lemma 5.**  $A_{\pi_F}$  defined in Lemma 2 is upper-bounded by a problem-dependent quantity,

$$A_{\pi_F} \leq \sum_{i \in \mathcal{K}} \left( \frac{32\alpha\sigma^2 \log_+(KT\delta_0^{-1/\alpha})}{\Delta_{i, h_{i,T}^+ - 1}} + \sqrt{32\alpha\sigma^2 \log_+(KT\delta_0^{-1/\alpha})} \right) + (K - 1)L.$$

*Proof.* We start from the result of Lemma 3,

$$A_{\pi_F} \leq \sum_{i \in \text{OP}_\xi} \left( 4 \sqrt{2\alpha\sigma^2 \log(KT\delta_0^{-1/\alpha})} \left( 1 + \sqrt{h_{i,T}^\xi - 1} \right) \right) + (K - 1)L. \quad (23)$$

We want to bound  $h_{i,T}^\xi$  with a problem dependent quantity  $h_{i,T}^+$ . We remind the reader that for arm  $i$  at round  $T$ , the  $h_{i,T}^\xi$ -th overpull has been on  $\xi_{t_i}$  pulled at round  $t_i$ . Therefore, Corollary 1 applies and we have

$$\begin{aligned} \bar{\mu}_{i, T}^{\xi, T-1} \left( N_{i,T}^{\pi^*} + h_{i,T}^\xi - 1 \right) &\geq \mu_T^+(\pi_F) - 4c \left( h_{i,T}^\xi - 1, \delta_{t_i} \right) \geq \mu_T^+(\pi_F) - 4c \left( h_{i,T}^\xi - 1, \delta_T \right) \\ &\geq \mu_T^+(\pi_F) - 4 \sqrt{\frac{2\alpha\sigma^2 \log(KT\delta_0^{-1/\alpha})}{h_{i,T}^\xi - 1}} \geq \mu_T^-(\pi^*) - 4 \sqrt{\frac{2\alpha\sigma^2 \log(KT\delta_0^{-1/\alpha})}{h_{i,T}^\xi - 1}}, \end{aligned}$$

with  $\mu_T^-(\pi^*) \triangleq \min_{i \in \mathcal{K}} \mu_i(N_{i,T}^{\pi^*} - 1)$  being the lowest mean reward for which a noisy value was ever obtained by the optimal policy.  $\mu_T^-(\pi^*) < \mu_T^+(\pi_F)$  implies that the regret is 0. Indeed, in that case the next possible pull with the largest mean for  $\pi_F$  is *strictly larger* than the mean of the last pull for  $\pi^*$ . Thus, there is no underpull at this round for  $\pi_F$  and  $R_T(\pi_F) = 0$  according to Equation 3. Therefore, we can assume  $\mu_T^-(\pi^*) \geq \mu_T^+(\pi_F)$  for the regret bound. Next, we define  $\Delta_{i,h} \triangleq \mu_T^-(\pi^*) - \bar{\mu}_i^h(N_{i,t}^{\pi^*} + h)$  as the difference between the lowest mean value of the arm pulled by  $\pi^*$  and the average of the  $h$  first overpulls of arm  $i$ . Thus, we have the following bound for  $h_{i,T}^\xi$ ,

$$h_{i,T}^\xi \leq 1 + \frac{32\alpha\sigma^2 \log(KT\delta_0^{-1/\alpha})}{\Delta_{i, h_{i,T}^\xi - 1}}.$$

Next,  $h_{i,T}^\xi$  has to be smaller than the maximum such  $h$ , for which the inequality just above is satisfied if we replace  $h_{i,T}^\xi$  by  $h$ . Therefore,

$$h_{i,T}^\xi \leq h_{i,T}^+ \triangleq \max \left\{ h \leq T \mid h \leq 1 + \frac{32\alpha\sigma^2 \log_+(KT\delta_0^{-1/\alpha})}{\Delta_{i,h-1}^2} \right\}. \quad (24)$$

Since the square-root function is increasing, we can upper-bound Equation 18 by replacing  $h_{i,T}^\xi$  by its upper bound  $h_{i,T}^+$  to get

$$\begin{aligned} A_{\pi_F} &\leq \sum_{i \in \text{OP}_\xi} \left( 4\sqrt{2\alpha\sigma^2 \log_+(KT\delta_0^{-1/\alpha})} \left( 1 + \sqrt{h_{i,T}^+ - 1} \right) + L \right) \\ &\leq \sum_{i \in \text{OP}_\xi} \left( \sqrt{32\alpha\sigma^2 \log_+(KT\delta_0^{-1/\alpha})} \left( 1 + \frac{\sqrt{32\alpha\sigma^2 \log_+(KT\delta_0^{-1/\alpha})}}{\Delta_{i,h_{i,T}^+-1}} \right) + L \right). \end{aligned}$$

The quantity  $\text{OP}_\xi$  is depends on the execution. Notice that there are at most  $K - 1$  arms in  $\text{OP}_\xi$  and that  $\text{OP} \subset \mathcal{K}$ . Therefore, we have

$$A_{\pi_F} \leq \sum_{i \in \mathcal{K}} \left( \frac{32\alpha\sigma^2 \log_+(KT\delta_0^{-1/\alpha})}{\Delta_{i,h_{i,T}^+-1}} + \sqrt{32\alpha\sigma^2 \log_+(KT\delta_0^{-1/\alpha})} \right) + (K - 1)L.$$

□

**Corollary 2.** For  $\delta_t \triangleq 1/(Kt^5)$  and  $C_\alpha \triangleq 32\alpha\sigma^2$ , the regret of FEWA is bounded as

$$\mathbb{E}[R_T(\pi_F)] \leq \sum_{i \in \mathcal{K}} \left( \frac{C_5 \log(KT)}{\Delta_{i,h_{i,T}^+-1}} + \sqrt{C_5 \log(KT)} + L \right).$$

*Proof.* Using Lemmas 2, 4, and 5 we get

$$\begin{aligned} \mathbb{E}[R_T(\pi_F)] &= \mathbb{E}[A_{\pi_F}] + \mathbb{E}[B] \leq \sum_{i \in \mathcal{K}} \left( \frac{32\alpha\sigma^2 \log(KT)}{\Delta_{i,h_{i,T}^+-1}} + \sqrt{32\alpha\sigma^2 \log(KT)} \right) + (K - 1)L + \frac{L\pi^2}{6} \\ &\leq \sum_{i \in \mathcal{K}} \left( \frac{32\alpha\sigma^2 \log(KT)}{\Delta_{i,h_{i,T}^+-1}} + \sqrt{32\alpha\sigma^2 \log(KT)} + L \right). \end{aligned}$$

□

**Corollary 4.** FEWA run with  $\alpha > 3$  and  $\delta_0 \triangleq 2\delta/\zeta(\alpha - 2)$  achieves with probability  $1 - \delta$ ,

$$R_T(\pi_F) \leq \sum_{i \in \mathcal{K}} \left( \frac{32\alpha\sigma^2 \log_+\left(\frac{KT\zeta(\alpha-2)^{1/\alpha}}{(2\delta)^{1/\alpha}}\right)}{\Delta_{i,h_{i,T}^+-1}} + \sqrt{32\alpha\sigma^2 \log_+\left(\frac{KT\zeta(\alpha-2)^{1/\alpha}}{(2\delta)^{1/\alpha}}\right)} \right) + (K - 1)L.$$

*Proof.* We consider the event  $\cup_{t \leq T} \xi_t$  which happens with probability

$$1 - \sum_{t \leq T} \frac{Kt^2\delta_t}{2} \leq 1 - \sum_{t \leq T} \frac{Kt^2\delta_t}{2} \leq 1 - \frac{\zeta(\alpha - 2)\delta_0}{2}.$$

Therefore, by setting  $\delta_0 \triangleq 2\delta/\zeta(\alpha - 2)$ , we have that with probability  $1 - \delta$ ,  $B = 0$  since  $[\xi_t^-] = 0$  for all  $t$ . We use Lemma 5 to get the claim of the corollary. □

## E Efficient algorithm EFF-FEWA

In Algorithm 3, we present EFF-FEWA, an algorithm that stores at most  $2K \log_2(t)$  statistics. More precisely, for  $j \leq \log_2(N_{i,t}^{\pi_{\text{EFF}}})$ , we let  $\widehat{s}_{i,j}^p$  and  $\widehat{s}_{i,j}^c$  be the current and pending  $j$ -th statistic for arm  $i$ . We then present an analysis of EFF-FEWA.

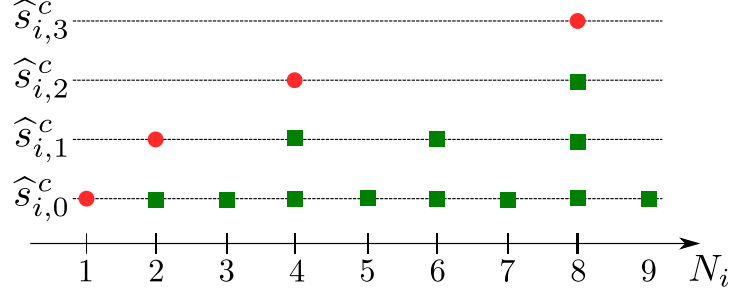


Figure 3: Illustration of the functioning of EFF-FEWA. The red circles denotes the number of pulls of arm  $i$  at which a new estimate  $\hat{s}_{i,j}^c$  is created corresponding to a window  $h = 2^j$ , while the green boxes indicate the number of pulls for which  $\hat{s}_{i,j}^c$  is updated with the last  $2^j$  samples.

---

**Algorithm 3** EFF-FEWA
 

---

**Input:**  $\mathcal{K}$ ,  $\delta_0$ ,  $\alpha$

- 1: pull each arm once, collect reward, and initialize  $N_{i,K} \leftarrow 1$
- 2: **for**  $t \leftarrow K + 1, K + 2, \dots$  **do**
- 3:    $\delta_t \leftarrow \delta_0 / (Kt^\alpha)$
- 4:    $j \leftarrow 0$  {initialize bandwidth}
- 5:    $\mathcal{K}_1 \leftarrow \mathcal{K}$  {initialize with all the arms}
- 6:    $i(t) \leftarrow \text{none}$
- 7:   **while**  $i(t)$  is **none** **do**
- 8:      $\mathcal{K}_{2^{j+1}} \leftarrow \text{EFF\_FILTER}(\mathcal{K}_{2^j}, j, \delta_t)$
- 9:      $j \leftarrow j + 1$
- 10:    **if**  $\exists i \in \mathcal{K}_{2^j}$  such that  $N_{i,t} \leq 2^j$  **then**
- 11:      $i(t) \leftarrow i$
- 12:    **end if**
- 13:   **end while**
- 14:   receive  $r_i(N_{i,t+1}) \leftarrow r_{i(t),t}$
- 15:   EFF\_UPDATE( $i(t), r_i(N_{i,t+1}), t + 1$ )
- 16: **end for**

---



---

**Algorithm 4** EFF\_FILTER
 

---

**Input:**  $\mathcal{K}_{2^j}$ ,  $j$ ,  $\delta_t$ ,  $\sigma$

- 1:  $c(2^j, \delta_t) \leftarrow \sqrt{2\sigma^2 / 2^j \log \delta_t^{-1}}$
- 2:  $\hat{s}_{\max,j}^c \leftarrow \max_{i \in \mathcal{K}_h} \hat{s}_{i,j}^c$
- 3: **for**  $i \in \mathcal{K}_h$  **do**
- 4:    $\Delta_i \leftarrow \hat{s}_{\max,j}^c - \hat{s}_{i,j}^c$
- 5:   **if**  $\Delta_i \leq 2c(2^j, \delta_t)$  **then**
- 6:     add  $i$  to  $\mathcal{K}_{2^{j+1}}$
- 7:   **end if**
- 8: **end for**

**Output:**  $\mathcal{K}_{2^{j+1}}$

---

---

**Algorithm 5** EFF\_UPDATE
 

---

**Input:**  $i, r, t$   
 1:  $N_{i(t),t} \leftarrow N_{i(t),t-1} + 1$   
 2:  $R_i^{\text{total}} \leftarrow R_i^{\text{total}} + r$  {keep track of total reward}  
 3: **if**  $\exists j$  such that  $N_{i,t} = 2^j$  **then**  
 4:    $\widehat{s}_{i,j}^c \leftarrow R_i^{\text{total}}/N_{i,t}$  {initialize new statistics}  
 5:    $\widehat{s}_{i,j}^p \leftarrow 0$   
 6:    $n_{i,j} \leftarrow 0$   
 7: **end if**  
 8: **for**  $j \leftarrow 0 \dots \log_2(N_{i,t})$  **do**  
 9:    $n_{i,j} \leftarrow n_i + 1$   
 10:    $\widehat{s}_{i,j}^p \leftarrow \widehat{s}_{i,j}^p + r$   
 11:   **if**  $n_{i,j} = 2^j$  **then**  
 12:      $\widehat{s}_{i,j}^c \leftarrow \widehat{s}_{i,j}^p/2^j$   
 13:      $n_{i,j} \leftarrow 0$   
 14:      $\widehat{s}_{i,j}^p \leftarrow 0$   
 15:   **end if**  
 16: **end for**

---

As  $N_{i,t}$  increases new statistics  $\widehat{s}_{i,j}^c$  for larger windows are created, as illustrated in Fig. 3. On one hand, at any time  $t$ ,  $\widehat{s}_{i,j}^c$  is the average of  $2^{j-1}$  consecutive reward samples for arm  $i$  within the last  $2^j - 1$  sample. These statistics are used in the filtering process as they are representative of exactly  $2^{j-1}$  recent samples. On the other hand,  $\widehat{s}_{i,j}^p$  stores the pending samples that are not yet taken into account by  $\widehat{s}_{i,j}^c$ . Therefore, each time we pull arm  $i$ , we update all the pending averages. When the pending statistic is the average of the  $2^{j-1}$  last samples then we set  $\widehat{s}_{i,j}^c \leftarrow \widehat{s}_{i,j}^p$  and we reinitialize  $\widehat{s}_{i,j}^p \leftarrow 0$ .

In analyzing the performance of EFF-FEWA, we have to account for two different effects: **1)** the loss in resolution due to windows of size that increases exponentially instead of fix increment of 1, and **2)** the delay in updating the statistics  $\widehat{s}_{i,j}^c$ , which do not include the most recent samples. We let  $\bar{\mu}_i^{h',h''}$  be the average of the samples between the  $h'$ -th to last one and the  $h''$ -th to last one (included) with  $h'' > h'$ . FEWA was controlling  $\bar{\mu}_i^{1,h}$  for each arm, EFF-FEWA controls  $\bar{\mu}_i^{h',h'+2^{j-1}}$  with different  $h'_i \leq 2^{j-1} - 1$  for each arm depending on when  $\widehat{s}_{i,j}^c$  was refreshed last time. However, since the means of arms are non-increasing, we can consider the worst case when the arm with the highest mean available at that round is estimated on its last samples (the smaller one) and the bad arms are estimated on their oldest possibles samples (the larger one).

**Lemma 6.** *On the favorable event  $\xi_t$ , if an arm  $i$  passes through a filter of window  $h$  at round  $t$ , the average of its  $h$  last pulls cannot deviate significantly from the best available arm  $i_t^*$  at that round,*

$$\bar{\mu}_i^{2^{j-1}, 2^j-1} \geq \mu_t^+(\pi_{\text{F}}) - 4c(h, \delta_t).$$

Then, we modify Corollary 1 to have the following efficient version of it.

**Corollary 5.** *Let  $i \in \text{OP}$  be an arm overpulled by EFF-FEWA at round  $t$  and  $h_{i,t}^{\pi_{\text{EF}}} \triangleq N_{i,t}^{\pi_{\text{EF}}} - N_{i,t}^{\pi^*} \geq 1$  be the difference in the number of pulls w.r.t. the optimal policy  $\pi^*$  at round  $t$ . On the favorable event  $\xi_t$ , we have that*

$$\mu_t^+(\pi_{\text{EF}}) - \bar{\mu}^{h_{i,t}^{\pi_{\text{EF}}}}(N_{i,t}) \leq \frac{4\sqrt{2}}{\sqrt{2}-1} c(h_{i,t}^{\pi_{\text{EF}}}, \delta_t).$$

*Proof.* If  $i$  was pulled at round  $t$ , then by the condition at Line 10 of Algorithm 3, it means that  $i$  passes through all the filters until at least window  $2^f$  such that  $2^f \leq h_{i,t}^{\pi_{\text{EF}}} < 2^{f+1}$ . Note that for  $h_{i,t}^{\pi_{\text{EF}}} = 1$ , then EFF-FEWA has



the same guarantee as FEWA since the first filter is always up to date. Then for  $h_{i,t}^{\pi_{\text{EF}}} \geq 2$ ,

$$\bar{\mu}_i^{1, h_{i,t}^{\pi_{\text{EF}}}}(N_{i,t}) \geq \bar{\mu}_i^{1, 2^f - 1}(N_{i,t}) = \frac{\sum_{j=1}^f 2^{j-1} \mu_i^{2^{j-1}, 2^j - 1}}{2^f - 1} \quad (25)$$

$$\geq \mu_t^+(\pi_{\text{EF}}) - \frac{4 \sum_{j=1}^f 2^{j-1} c(2^{j-1}, \delta)}{2^f - 1} = \mu_t^+(\pi_{\text{EF}}) - 4c(1, \delta_t) \frac{\sum_{j=1}^f \sqrt{2}^{j-1}}{2^f - 1} \quad (26)$$

$$= \mu_t^+(\pi_{\text{EF}}) - 4c(1, \delta_t) \frac{\sqrt{2}^f - 1}{(2^f - 1)(\sqrt{2} - 1)} \geq \mu_t^+(\pi_{\text{EF}}) - 4c(1, \delta_t) \frac{1}{\sqrt{2}^f (\sqrt{2} - 1)} \quad (27)$$

$$= \mu_t^+(\pi_{\text{EF}}) - \frac{4\sqrt{2}}{\sqrt{2} - 1} c(2^{f+1}, \delta_t) \geq \mu_t^+(\pi_{\text{EF}}) - \frac{4\sqrt{2}}{\sqrt{2} - 1} c(h_{i,t}^{\pi_{\text{EF}}}, \delta_t), \quad (28)$$

where Equation 25 uses that the average of older means is larger than average of the more recent ones and then decomposes  $2^f - 1$  means onto a geometric grid. Then, Equation 26 uses Lemma 6 and make the dependence of  $c(2^{j-1}, \delta)$  on  $j$  explicit. Next, Equations 27 and 28 use standard algebra to derive a lower bound and that  $c(h, \delta)$  decreases with  $h$ .  $\square$

Using the result above, we follow the same proof as the one for FEWA and derive minimax and problem-dependent upper bounds for EFF-FEWA using Corollary 5 instead of Corollary 1.

**Corollary 6** (minimax guarantee for EFF-FEWA). *For any rotting bandit scenario with means  $\{\mu_i(n)\}_{i,n}$  satisfying Assumption 1 with bounded decay  $L$  and any time horizon  $T$ , EFF-FEWA with  $\delta_t = 1/(Kt^5)$ ,  $\alpha = 5$ , and  $\delta_0 = 1$ , has its expected regret upper-bounded as*

$$\mathbb{E}[R_T(\pi_{\text{EF}})] \leq 13\sigma \left( \frac{\sqrt{2}}{\sqrt{2} - 1} \sqrt{KT} + K \right) \sqrt{\log(KT)} + KL.$$

**Corollary 7** (problem-dependent guarantee for EFF-FEWA). *For  $\delta_t = 1/(Kt^5)$ , the regret of EFF-FEWA is upper-bounded as*

$$R_T(\pi_{\text{EF}}) \leq \sum_{i \in \mathcal{K}} \left( \frac{C_5 \frac{2}{3-2\sqrt{2}} \log(KT)}{\Delta_{i, h_{i,T}^+}} + \sqrt{C_5 \log(KT)} + L \right),$$

with  $C_\alpha \triangleq 32\alpha\sigma^2$  and  $h_{i,T}^+$  defined in Equation 10.

## F Numerical simulations: Stochastic bandits

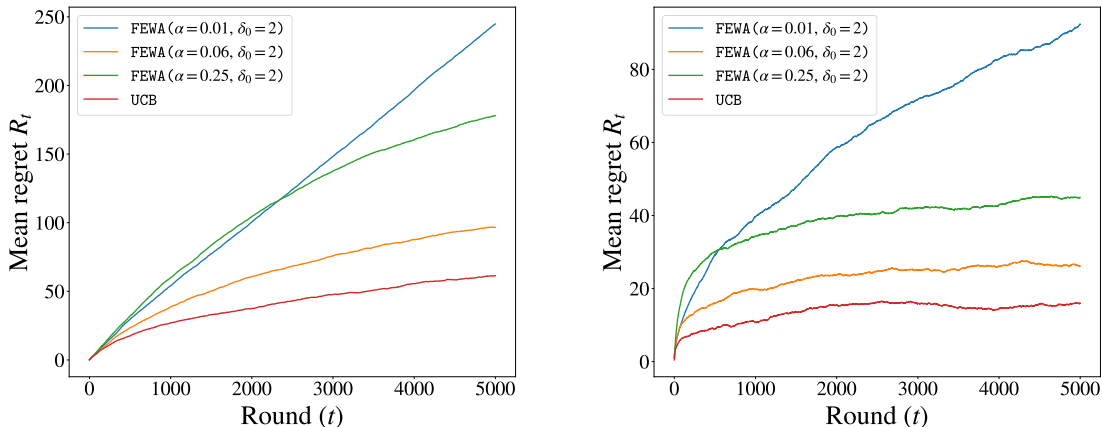


Figure 4: Comparing UCB1 and FEWA with  $\Delta = 0.14$  and  $\Delta = 1$ .

In Figure 4 we compare the performance of FEWA against UCB1 (Auer et al., 2002a) on two-arm bandits with different gaps. These experiments confirm the theoretical findings of Theorem 1 and Corollary 2: FEWA has comparable performance with UCB1. In particular, both algorithms have a logarithmic asymptotic behavior and for  $\alpha = 0.06$ , the ratio between the regret of two algorithms is empirically lower than 2. Notice, the theoretical factor between the two upper bounds is 5 (for  $\alpha = 5$ ). This shows the ability of FEWA to be competitive for stochastic bandits.