
Regret Bounds and Reinforcement Learning Exploration of EXP-based Algorithms

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 EXP-based algorithms are often used for exploration in non-stochastic bandit
2 problems assuming rewards are bounded. Motivated by the recent advancements
3 in reinforcement learning with rewards of any scale, we propose a new algorithm,
4 namely EXP4.P, by modifying EXP4 and establish its regret upper bounds in
5 both bounded and unbounded sub-Gaussian contextual bandits. The unbounded
6 reward result also holds for a revised version of EXP3.P. Moreover, we provide
7 a lower bound on regret that suggests no sublinear regret can be achieved given
8 short time horizon. Unbounded rewards pose challenges as the regret cannot
9 be limited by the number of trials, and choosing suboptimal arms may result in
10 infinite regret. We also extend EXP4.P from bandit to reinforcement learning to
11 incentivize exploration by multiple agents given black-box rewards. The resulting
12 algorithm has been tested on hard-to-explore games and it shows an improvement
13 on exploration compared to state-of-the-art.

14 1 Introduction

15 Multi-armed bandit (MAB) is to maximize cumulative reward of a player throughout a bandit game
16 by choosing different arms at each time step. It is also equivalent to minimizing the regret defined
17 as the difference between the best rewards that can be achieved and the actual reward gained by
18 the player. Formally, given time horizon T , in time step $t \leq T$ the player chooses one arm a_t
19 among K arms, receives $r_{a_t}^t$ among rewards $r^t = (r_1^t, r_2^t, \dots, r_K^t)$, and maximizes the total reward
20 $\sum_{t=1}^T r_{a_t}^t$ or minimizes the regret. Traditionally, there are two classic versions of non-stochastic
21 bandits: Adversarial and Contextual. For adversarial MAB, rewards of the K arms r^t can be chosen
22 arbitrarily by adversaries at step t . When the adversary is a context-dependent reward generator, it
23 boils down to contextual bandits. Contextual bandit is a variant of MAB by adding context or state
24 space S and a different regret definition. At time step t , the player has context $s_t \in S$ and rewards r^t
25 follow $f(\mu(s_t))$ where f is any distribution and $\mu(s_t)$ is the mean vector that depends on state s_t .

26 Computationally efficient and with abundant theoretical analyses are the EXP-type MAB algorithms.
27 Specifically, the regret of EXP3.P for adversarial bandit achieves optimality both in the expected and
28 high probability sense. In EXP3.P, each arm has a trust coefficient (weight). The player samples each
29 arm with probability being the sum of its normalized weights and a bias term, receives reward of the
30 sampled arm and exponentially updates the weights based on the corresponding reward estimates.
31 It achieves the regret of the order $O(\sqrt{T})$ in a high probability sense. To incorporate the context
32 information in contextual bandits, a variant of EXP-type algorithms is proposed as EXP4 [3]. In
33 EXP4, there are any number of experts. Each expert has a sample rule over actions (arms) and a
34 trust coefficient. The player samples according to the weighted average of experts' sample rules and
35 updates the weights respectively. Then the regret is defined by comparing the actual reward with the
36 reward that can be achieved by the best expert instead of by the best arm. The expectation of regret is
37 proven to be optimal for contextual bandit. Independently, [11] propose a modification of EXP4 that

38 achieves high probability guarantee, which, however, requires changes in the reward estimates. A
39 high probability regret has not yet been studied in its original form of EXP4.

40 Recently, contextual bandit has been further aligned with Reinforcement Learning (RL) where state
41 and reward transitions follow a Markov Decision Process (MDP) represented by transition kernel
42 $P(s_{t+1}, r^t | a_t, s_t)$. A key challenge in RL is the trade-off between exploration and exploitation.
43 Exploration is to encourage the player to try new arms in bandit or new actions in RL to understand
44 the game better. It helps to plan for the future, but with the sacrifice of potentially lowering the current
45 reward. Exploitation aims to exploit currently known states and arms to maximize the current reward,
46 but it potentially prevents the player to gain more information to increase future reward. To maximize
47 the cumulative reward, the player needs to learn the game by exploration, while guaranteeing current
48 reward by exploitation.

49 How to incentivize exploration in RL has been a main focus in RL. Since RL is built on bandits, it
50 is natural to extend bandit techniques to RL and UCB is such a success. UCB [2] motivates count-
51 based exploration [18] in RL and the subsequent Pseudo-Count exploration [4], though it is initially
52 developed for stochastic bandits. Another line of work on RL exploration is based on deep learning
53 techniques. Using deep neural networks to keep track of the Q -values by means of Q -networks in RL
54 is called DQN [9]. This combination of deep learning and RL has shown great success. ϵ -greedy in
55 [10] is a simple exploration technique based on DQN. Besides ϵ -greedy, intrinsic model exploration
56 computes intrinsic rewards that directly measure and thereby incentivizing exploration when added to
57 extrinsic (actual) rewards of RL, e.g. DORA [6] and [17]. Random Network Distillation (RND) [5]
58 is a more recent suggestion relying on a fixed target network. A drawback of RND is its local focus
59 without global exploration. EXP-type algorithms in contextual bandits work by integrating arbitrary
60 experts and hence providing exploration possibilities for RL, which, however, has not yet been studied.
61 Furthermore, the existing EXP4 or its variant cannot be directly adapted to RL. It is worth noting
62 that EXP-type algorithms are optimal under the assumption that $0 \leq r_i^t \leq 1$ for any arm i and step t .
63 The uniformly bounded assumption is crucial in the proof of regret bounds for existing EXP-type
64 algorithms. It requires the rewards to be scalable with the knowledge of a uniform bound for all
65 rewards in all states or context vectors. Nevertheless, reward in RL can be unbounded and unscalable
66 in real-world scenarios, which violates the bounded assumption. Examples include navigation tasks,
67 where the reward is unbounded for each step that brings the agent closer to the goal, and racing tasks,
68 where the reward is the distance covered by the agent. The counterpart of bandit algorithms in the
69 unbounded or scale-free case remained unexplored, unite the work herein and it necessitates a new
70 algorithm based on EXP3.P and EXP4.

71 In this paper, we are the first to propose a new algorithm, EXP4.P based on EXP4 without changing
72 the reward estimates. We show its optimal regret holds with high probability and in expectation for
73 contextual bandits with possibly unbounded (scale-free) rewards. The regret bounds for unbounded
74 bandits studied herein are significantly different from prior works. Compared to the high probability
75 version in [11], our algorithm only requires one parameter, is consistent with the reward estimate in
76 EXP4 and EXP3.P, removes the reward assumption of $[0, 1]$, and generalizes to the expected regret.
77 The proof extension to the unbounded case is non-trivial since it requires several deep results from
78 information theory and probability, by first establishing a high probability regret in the bounded case
79 with exponential terms and then using the Randemacher complexity theory to capture the dynamics
80 of arm selection in the unbounded case. Combining all these together is very technical and requires
81 new ideas. As a by-product, the analysis can be applied to EXP3.P to deliver a similar result for
82 bandits without expert advice. The upper bound for unbounded bandits requires T to be sufficiently
83 large, i.e. unbounded rewards may lead to extremely large regret without enough exploration, which
84 is computationally expensive in an RL setting. We herein provide a worst-case analysis implying no
85 sublinear regret can be achieved below an instance-specific minimal T , by our brand new construction
86 of instances. Precisely, we derive lower bounds of order $\Omega(T)$ for certain fixed T and upper bounds
87 of order $O^*(\sqrt{T})$ for T being large enough. The question of bounds for any value of T remains open.

88 Given the challenges of RL context where rewards are possibly unbounded or unrescalable which have
89 not been addressed by existing methods, we combine the proposed scale-free EXP-type algorithms
90 with deep RL. To this end, we extend the new EXP4.P to RL that allows for general experts by
91 generalizing the concept of experts to be any RL algorithms. Here experts improve local policies
92 with the underlying Markov process and exponential weights are assigned to the experts to produce a
93 global optimal policy. This is the first RL algorithm using several experts enabling global exploration,
94 where the overall performance is comparable to the best model even if we do not know which one

95 is the best beforehand, and thereby achieving model selections [8]. To address the issue of EXP4’s
 96 inefficiency with a large number of experts, we combine EXP4-RL with at least one state-of-the-art
 97 expert algorithm for improved efficiency and performance thus having only a few experts. Focusing
 98 on DQN, in the computational study we focus on two agents consisting of RND and ϵ -greedy DQN.
 99 We implement the EXP4-RL algorithm on hard-to-explore RL games Montezuma’s Revenge and
 100 Mountain Car and compare it with the benchmark RND [5]. The numerical results show that the
 101 algorithm gains more exploration than RND and it gains the ability of global exploration by avoiding
 102 local maxima of RND. Its total reward also increases with training. Overall, our algorithm improves
 103 exploration on the benchmark games.

104 The main contributions are as follows. We introduce sub-Gaussian bandits with the unique aspect and
 105 challenge of unbounded and scale-free rewards both in contextual bandits and MAB when EXP-based
 106 algorithms are considered. We propose a new EXP4.P algorithm based on EXP4 and EXP3.P and
 107 analytically establish its optimal regret both in unbounded and bounded cases. Unbounded rewards
 108 and contextual setting pose non-trivial challenges in the analyses. We also provide the very first regret
 109 lower bound in such a case that indicates a threshold of T for sublinear regret, by constructing a novel
 110 family of Gaussian bandits. We also provide the very first extension of EXP4.P to RL exploration
 111 using multiple agents and show its superior performance on two hard-to-explore RL games.

112 A literature review is provided in Section 2. Then in Section 3 we develop a new algorithm EXP4.P
 113 by modifying EXP4, and exhibit its regret bounds for contextual bandits and that of the EXP3.P
 114 algorithm for unbounded MAB, and lower bounds. Section 4 discusses the EXP4.P algorithm for RL
 115 exploration. Finally, in Section 5, we present numerical results related to the proposed algorithm.

116 2 Literature Review

117 The importance of exploration in RL is well understood. Count-based exploration in RL is such
 118 a success with the UCB technique. [18] develop Bellman value iteration $V(s) = \max_a \hat{R}(s, a) +$
 119 $\gamma E[V(s')] + \beta N(s, a)^{-\frac{1}{2}}$, where $N(s, a)$ is the number of visits to (s, a) for state s and action
 120 a . Value $N(s, a)^{-\frac{1}{2}}$ is positively correlated with curiosity of (s, a) and encourages exploration.
 121 This method is limited to tableau model-based MDP for small state spaces. While [4] introduce
 122 Pseudo-Count exploration for non-tableau MDP with density models, it is hard to model. However,
 123 UCB achieves optimality if bandits are stochastic and may suffer linear regret otherwise [21]. The
 124 work on CORRAL in [1] considers a group of bandit algorithms, but it requires a parameter search in
 125 the parameter space. In the RL setting, such updates are inefficient and do not fit the dynamic RL
 126 setting. EXP-type algorithms for non-stochastic bandits can generalize to RL with fewer assumptions
 127 about the statistics of rewards, which have not yet been studied. In conjunction with DQN, ϵ -greedy
 128 in [10] is a simple exploration technique using DQN. Besides ϵ -greedy, intrinsic model exploration
 129 computes intrinsic rewards by the accuracy of a model trained on experiences. Intrinsic rewards
 130 directly measure and incentivize exploration if added to actual rewards of RL, e.g. see [6, 17, 5].
 131 Random Network Distillation(RND) in [5] define it as $e(s', a) = \|\hat{f}(s') - f(s')\|_2^2$ where \hat{f} is a
 132 parametric model and f is a randomly initialized but fixed model. Here $e(s', a)$, independent of the
 133 transition, only depends on state s' and drives RND to outperform others on Montezuma’s Revenge.
 134 None of these algorithms use several experts which is a significant departure from our work.

135 Along the line of work on regret analyses focusing on EXP-type algorithms, [3] first introduce
 136 EXP3.P for bounded adversarial MAB and EXP4 for bounded contextual bandits. For the EXP3.P
 137 algorithm, an upper bound on regret of order $O(\sqrt{T})$ holds with high probability and in expectation,
 138 which has no gap with the lower bound and hence it establishes that EXP3.P is optimal. EXP4 is
 139 optimal for contextual bandits in the sense that its expected regret is $O(\sqrt{T})$. Then [11] extend
 140 it to a high probability counterpart by modifying the reward estimates. These regret bounds are
 141 invalid for bandits with unbounded support. Though [16] demonstrate a regret bound $O(\sqrt{T \cdot \gamma_T})$
 142 for noisy Gaussian process bandits, information gain γ_T is not well-defined in a noiseless setting.
 143 For noiseless Gaussian bandits, [7] show both the optimal lower and upper bounds on regret, but
 144 the regret definition is not consistent with [3]. We tackle these problems by establishing an upper
 145 bound of order $O^*(\sqrt{T})$ on regret 1) with high probability for bounded contextual bandit and 2) for
 146 sub-Gaussian bandit both in expectation and with high probability.

147 3 Regret Bounds

148 We first introduce notations. Let T be the time horizon. For bounded bandits, at step t , $0 < t \leq T$
 149 rewards r^t can be chosen arbitrarily under the condition that $-1 \leq r^t \leq 1$. For unbounded bandits,
 150 let rewards r^t follow multi-variate distribution $f_i(\mu, \Sigma)$ where $\mu = (\mu_1, \mu_2, \dots, \mu_K)$ is the mean
 151 vector and $\Sigma = (a_{ij})_{i,j \in \{1, \dots, K\}}$ is the covariance matrix of the K arms and f_i is the density.
 152 We specify f_t to be non-degenerate sub-Gaussian for analyses on light-tailed distributions where
 153 $\min_j a_{j,j} > 0$. A random variable X is σ^2 -sub-Gaussian if for any $t > 0$, the tail probability satisfies
 154 $P(|X| > t) \leq B e^{-\sigma^2 t^2}$ where B is a positive constant.

155 The player receives reward $y_t = r_{a_t}^t$ by pulling arm a_t . The regret is defined as $R_T =$
 156 $\max_j \sum_{t=1}^T r_j^t - \sum_{t=1}^T y_t$ in adversarial bandits that depends on realizations of rewards. For con-
 157 textual bandits with experts, besides the above let N be the number of experts and c_t be the context
 158 information. We denote the reward of expert i by $G_i = \sum_{t=1}^T z_i(t) = \sum_{t=1}^T \xi_i(t)^T x(t)$, where
 159 $x(t) = r^t$ and $\xi_i(t) = (\xi_i^1(t), \dots, \xi_i^K(t))$ is the probability vector of expert i . Then regret is defined
 160 as $R_T = \max_i G_i - \sum_{t=1}^T y_t$, which is with respect to the best expert, rather than the best arm in
 161 MAB. This is reasonable since a uniform optimal arm is a special expert assigning probability 1 to the
 162 optimal arm throughout the game and experts can potentially perform better and admit higher rewards.
 163 This coincides with our generalization of EXP4.P to RL where the experts can be well-trained neural
 164 networks. We follow established definitions of pseudo regret $R'_T = T \cdot \max_k \mu_k - \sum_t E[y_t]$ and
 165 $\sum_{t=1}^T \max_i \sum_{j=1}^K \xi_i^j(t) \mu_j - \sum_t E[y_t]$ in adversarial and contextual bandits, respectively.

166 3.1 Contextual Bandits and EXP4.P Algorithm

167 For contextual bandits, [3] give the EXP4 algorithm and prove its expected regret to be optimal
 168 under the bounded assumption on rewards and under the assumption that a uniform expert is always
 169 included, where by uniform expert we refer to an expert that always assigns equal probability to each
 170 arm. Our goal is to extend EXP4 to RL where rewards are often unbounded, such as several games
 171 in OpenAI gym, for which the theoretical guarantee of EXP4 may be absent. To this end, herein
 172 we propose a new Algorithm, named EXP4.P, as a variant of EXP4. Its effectiveness is two-fold.
 173 First, we show that EXP4.P has an optimal regret with high probability in the bounded case and
 174 consequently, we claim that the regret of EXP4.P is still optimal given unbounded bandits. All the
 175 proof are in the Appendix under the aforementioned assumption on experts. Second, it is successfully
 176 extended to RL where it achieves computational improvements.

177 3.1.1 EXP4.P Algorithm

Algorithm 1 EXP4.P

Initialization: Weights $w_i(1) = \exp(\frac{\alpha \gamma}{3K} \sqrt{NT})$, $i \in \{1, 2, \dots, N\}$ for $\alpha > 0$ and $\gamma \in (0, 1)$;
for $t = 1, 2, \dots, T$ **do**

 Get probability vectors $\xi_1(t), \dots, \xi_N(t)$ of arms from experts where $\xi_i(t) = (\xi_i^j(t))_j$;

 For any $j = 1, 2, \dots, K$, set $p_j(t) = (1 - \gamma) \sum_{i=1}^N \frac{w_i(t) \cdot \xi_i^j(t)}{\sum_{j=1}^N w_j(t)} + \frac{\gamma}{K}$;

 Choose i_t randomly according to the distribution $p_1(t), \dots, p_K(t)$;

 Receive reward $r_{i_t}(t) = x_{i_t}(t)$;

 For any $j = 1, \dots, K$, set $\hat{x}_j(t) = \frac{r_j(t)}{p_j(t)} \cdot \mathbf{1}_{j=i_t}$;

 Set $\hat{x}(t) = (\hat{x}_j(t))_j$;

 For any $i = 1, \dots, N$, set

$$\hat{z}_i(t) = \xi_i(t)^T \hat{x}(t) \text{ and } w_i(t+1) = w_i(t) \exp(\frac{\gamma}{3K} (\hat{z}_i(t) + \frac{\alpha}{(\sum_{j=1}^N w_j(t) + \frac{\gamma}{K}) \sqrt{NT}}));$$

end for

178 Our proposed EXP4.P is shown as Algorithm 1. The main modifications compared to EXP4 lie in the
 179 update and the initialization of trust coefficients of experts as highlighted. The upper bound of the
 180 confidence interval of the reward estimate is added to the update rule for each expert, in the spirit of
 181 EXP3.P (see Algorithm 2) and removing the need of changing the reward estimate. However, this
 182 term and initialization of EXP4.P are quite different from that in EXP3.P for MAB.

183 **3.1.2 Bounded Rewards**

184 Borrowing the ideas of [3], we claim EXP4.P has an optimal sublinear regret with high probability
 185 by first establishing two lemmas presented in Appendix. The main theorem is as follows. We assume
 186 that the expert family includes a uniform expert, which is also assumed in the analysis of EXP4 in [3].

187 **Theorem 1.** *Let $0 \leq r^t \leq 1$ for every t . For any fixed time horizon $T > 0$, for all $K, N \geq 2$ and
 188 for any $1 > \delta > 0$, $\gamma = \sqrt{\frac{3K \ln N}{T(\frac{2N}{3} + 1)}} \leq \frac{1}{2}$, $\alpha = 2\sqrt{K \ln \frac{NT}{\delta}}$, we have that with probability at least
 189 $1 - \delta$, $R_T \leq 2\sqrt{3KT\left(\frac{2N}{3} + 1\right) \ln N} + 4K\sqrt{KNT \ln\left(\frac{NT}{\delta}\right)} + 8NK \ln\left(\frac{NT}{\delta}\right)$.*

190 Theorem 1 implies $R_T \leq O^*(\sqrt{T})$. The regret bound does depend on N . In practice the number of
 191 experts is small compared to the time horizon and the independence among experts makes parallelism
 192 a possibility. Note that $\gamma < \frac{1}{2}$ for large enough T . The proof of Theorem 1 essentially relies on the
 193 convergence of the reward estimators, similar to that in [3]. However, the objectives are different
 194 from [3], since our estimations and update of trust coefficients in EXP4.P are for experts, instead of
 195 EXP3.P for arms. This characterizes the relationships among EXP4.P estimates and the actual value
 196 of experts' rewards and the total rewards gained by EXP4.P and brings non-trivial challenges.

197 **3.1.3 Unbounded Rewards**

198 We proceed to show optimal regret bounds of EXP4.P for unbounded contextual bandit. Again, a
 199 uniform expert is assumed to be included in the expert family. Surprisingly, we report that the analysis
 200 can be adapted to the existing EXP3.P in next section, which leads to optimal regret in MAB under
 201 no bounded assumption which is also a new result.

202 **Theorem 2.** *For sub-Gaussian bandits, any time horizon T , for any $0 < \eta < 1$, $0 < \delta < 1$
 203 and γ, α as in Theorem 1, with probability at least $(1 - \delta)(1 - \eta)^T$, EXP4.P has regret $R_T \leq$
 204 $4\Delta(\eta)\left(2\sqrt{3KT\left(\frac{2N}{3} + 1\right) \ln N}\right) + 4\Delta(\eta)\left(4K\sqrt{KNT \ln\left(\frac{NT}{\delta}\right)} + 8NK \ln\left(\frac{NT}{\delta}\right)\right)$ where $\Delta(\eta)$
 205 is determined by $\int_{-\Delta}^{\Delta} \dots \int_{-\Delta}^{\Delta} f(x_1, \dots, x_K) dx_1 \dots dx_K = 1 - \eta$ which yields $\Delta(\eta)$ of $O\left(\frac{1}{a} \log \frac{1}{\eta}\right)$.*

206 In the proof of Theorem 2, we first perform truncation of the rewards of sub-Gaussian bandits by
 207 dividing the rewards to a bounded part and unbounded tail. For the bounded part, we directly apply
 208 the upper bound on regret of EXP4.P presented in Theorem 1 and conclude with the regret upper
 209 bound of order $O(\Delta(\eta)\sqrt{T})$. Since a sub-Gaussian distribution is a light-tailed distribution we can
 210 control the probability of the tail, i.e. the unbounded part, which leads to the overall result.

211 The dependence of the bound on Δ can be removed by considering large enough T as stated next.

212 **Theorem 3.** *For sub-Gaussian bandits, for any $a > 2$, $0 < \delta < 1$, and γ, α as in Theorem 1, EXP4.P
 213 has regret $R_T \leq \log(1/\delta)O^*(\sqrt{T})$ with probability $(1 - \delta) \cdot \left(1 - \frac{1}{T^a}\right)^T$.*

214 Note that the constant term in $O^*(\cdot)$ depends on a . The above theorems deal with R_T ; an upper
 215 bound on pseudo regret or expected regret is established next. It is easy to verify by the Jensen's
 216 inequality that $R'_T \leq E[R_T]$ and thus it suffices to obtain an upper bound on $E[R_T]$.

217 For bounded bandits, the upper bound for $E[R_T]$ is of the same order as R_T which follows by a
 218 simple argument. For sub-Gaussian bandits, establishing an upper bound on $E[R_T]$ or R'_T based
 219 on R_T requires more work. We show an upper bound on $E[R_T]$ by using certain inequalities, limit
 220 theories, and Rademacher complexity. To this end, the main result reads as follows.

221 **Theorem 4.** *The regret of EXP4.P for sub-Gaussian bandits satisfies $R'_T \leq E[R_T] \leq O^*(\sqrt{T})$
 222 under the assumptions stated in Theorem 3.*

223 **3.2 MAB and EXP3.P Algorithm**

224 In this section, we establish upper bounds on regret in MAB given a high probability regret bound
 225 achieved by EXP3.P in [3]. We revisit EXP3.P and analyze its regret in unbounded scenarios in line
 226 with EXP4.P. Formally, we show that EXP3.P achieves regret of order $O^*(\sqrt{T})$ in sub-Gaussian
 227 MAB, with respect to R_T , $E[R_T]$ and R'_T . The results are summarized as follows.

228 **Theorem 5.** For sub-Gaussian MAB, any T , for any $0 < \eta, \delta < 1$, $\gamma = 2\sqrt{\frac{3K \ln K}{5T}}$, $\alpha = 2\sqrt{\ln \frac{NT}{\delta}}$,
229 EXP3.P has regret $R_T \leq 4\Delta(\eta) \cdot (\sqrt{KT \log(\frac{KT}{\delta})} + 4\sqrt{\frac{5}{3}KT \log K} + 8 \log(\frac{KT}{\delta}))$ with probability
230 $(1 - \delta)(1 - \eta)^T$ where $\Delta(\eta) = O(\frac{1}{a} \log \frac{1}{\eta})$, i.e. $\int_{-\Delta}^{\Delta} \dots \int_{-\Delta}^{\Delta} f(x_1, \dots, x_K) dx_1 \dots dx_K = 1 - \eta$.
231 To proof Theorem 5, we again do truncation. We apply the bounded result of EXP3.P in [3] and
232 achieve a regret upper bound of order $O(\Delta(\eta)\sqrt{T})$. The proof is similar to the proof of Theorem 2
233 for EXP4.P.
234 Similarly, we remove the dependence of the bound on Δ in Theorem 6 and claim a bound on the
expected regret for sufficiently large T in Theorem 7.

Algorithm 2 EXP3.P

Initialization: Weights $w_i(1) = \exp(\frac{\alpha\gamma}{3}\sqrt{\frac{T}{K}})$, $i \in \{1, 2, \dots, K\}$ for $\alpha > 0$ and $\gamma \in (0, 1)$;
for $t = 1, 2, \dots, T$ **do**
For any $i = 1, 2, \dots, K$, set $p_i(t) = (1 - \gamma) \frac{w_i(t)}{\sum_{j=1}^K w_j(t)} + \frac{\gamma}{K}$;
Choose i_t randomly according to the distribution $p_1(t), \dots, p_K(t)$;
Receive reward $r_{i_t}(t)$;
For $1 \leq j \leq K$, set $\hat{x}_j(t) = \frac{r_j(t)}{p_j(t)} \cdot \mathbb{1}_{j=i_t}$ and $w_j(t+1) = w_j(t) \exp(\frac{\gamma}{3K}(\hat{x}_j(t) + \frac{\alpha}{p_j(t)\sqrt{KT}}))$;
end for

235 **Theorem 6.** For sub-Gaussian MAB, for $a > 2$, $0 < \delta < 1$, and γ, α as in Theorem 5, EXP3.P has
236 regret $R_T \leq \log(1/\delta)O^*(\sqrt{T})$ with probability $(1 - \delta) \cdot (1 - \frac{1}{T^a})^T$.

238 **Theorem 7.** The regret of EXP3.P in sub-Gaussian MAB satisfies $R'_T \leq E[R_T] \leq O^*(\sqrt{T})$ with
239 the same assumptions as in Theorem 6.

240 3.3 Lower Bounds on Regret

241 Algorithms can suffer extremely large regret without enough exploration when playing unbounded
242 bandits given small T . To argue that our bounds on regret are not loose, we derive a lower bound on
243 the regret for sub-Gaussian bandits that essentially suggests that no sublinear regret can be achieved
244 if T is less than an instance-dependent bound. The main technique is to construct instances that have
245 certain regret, no matter what strategies are deployed. We need the following assumption.

246 **Assumption 1** There are two types of arms with general K with one type being superior (S is
247 the set of superior arms) and the other being inferior (I is the set of inferior arms). Let $1 - q, q$
248 be the proportions of the superior and inferior arms, respectively which is known to the adversary
249 and clearly $0 \leq q \leq 1$. The arms in S are indistinguishable and so are those in I . The first pull
250 of the player has two steps. First the player selects an inferior or superior set of arms based on
251 $P(S) = 1 - q, P(I) = q$ and once a set is selected, the corresponding reward of an arm from the
252 selected set is received.

253 An interesting special case of Assumption 1 is the case of two arms and $q = 1/2$. In this case, the
254 player has no prior knowledge and in the first pull chooses an arm uniformly at random.

255 The lower bound is defined as $R_L(T) = \inf \sup R'_T$, where, first, \inf is taken among all the strategies
256 and then \sup is among all Gaussian MAB. The following is the main result for lower bounds based
257 on inferior arms being distributed as $\mathcal{N}(0, 1)$ and superior as $\mathcal{N}(\mu, 1)$ with $\mu > 0$.

258 **Theorem 8.** In Gaussian MAB under Assumption 1, for any $q \geq 1/3$ we have $R_L(T) \geq (q -$
259 $\epsilon) \cdot \mu \cdot T$, where μ has to satisfy $G(q, \mu) < q$ with ϵ and T determined by $G(q, \mu) < \epsilon < q, T \leq$
260 $\frac{\epsilon - G(q, \mu)}{(1-q) \cdot \int |e^{-\frac{x^2}{2}} - e^{-\frac{(x-\mu)^2}{2}}| + 2 \text{ where } G(q, \mu) \text{ is } \max\{\int |qe^{-\frac{x^2}{2}} - (1-q)e^{-\frac{(x-\mu)^2}{2}}| dx,$
261 $\int |(1-q)e^{-\frac{x^2}{2}} - qe^{-\frac{(x-\mu)^2}{2}}| dx\}$.

262 To prove Theorem 8, we construct a special subset of Gaussian MAB with equal variances and zero
263 covariances. On these instances we find a unique way to explicitly represent any policy. This builds a
264 connection between abstract policies and this concrete mathematical representation. Then we show

265 that pseudo regret R'_T must be greater than certain values no matter what policies are deployed, which
 266 indicates a regret lower bound on this subset of instances.

267 Feasibility of the aforementioned conditions is established in the following theorem.

268 **Theorem 9.** *In Gaussian MAB under Assumption 1, for any $q \geq 1/3$, there exist μ and $\epsilon, \epsilon < \mu$ such*
 269 *that $R_L(T) \geq (q - \epsilon) \cdot \mu \cdot T$.*

270 The following result with two arms and equal probability in the first pull deals with general MAB. It
 271 shows that for any fixed $\mu > 0$ there is a minimum T and instances of MAB so that no algorithm can
 272 achieve sublinear regret. Table 1 (see Appendix) exhibits how the threshold of T varies with μ .

273 **Theorem 10.** *For general MAB under Assumption 1 with $K = 2, q = 1/2$, we have that $R_L(T) \geq$*
 274 *$\frac{T \cdot \mu}{4}$ holds for any distributions f_0 for the arms in I and f_1 for the arms in S with $\int |f_1 - f_0| > 0$*
 275 *(possibly with unbounded support), for any $\mu > 0$ and T satisfying $T \leq \frac{1}{2 \cdot \int |f_0 - f_1|} + 1$.*

276 4 EXP4.P Algorithm for RL

277 EXP4 has shown effectiveness in contextual bandits with statistical validity. Therefore, in this section,
 278 we extend EXP4.P to RL in Algorithm 3 where rewards are assumed to be nonnegative.

279 The player has experts that are represented by deep Q -networks trained by RL algorithms (there
 280 is a one to one correspondence between the experts and Q -networks). Each expert also has a trust
 281 coefficient. Trust coefficients are also updated exponentially based on the reward estimates as in
 282 EXP4.P. At each step of one episode, the player samples an expert (Q -network) with probability that
 283 is proportional to the weighted average of expert's trust coefficients. Then ϵ -greedy DQN is applied
 284 on the chosen Q -network. Here different from EXP4.P, the player needs to store all the interaction
 285 tuples in the experience buffer since RL is a MDP. After one episode, the player trains all Q -networks
 with the experience buffer and uses the trained networks as experts for the next episode. The basic

Algorithm 3 EXP4-RL

Initialization: Trust coefficients $w_k = 1$ for any $k \in \{1, \dots, E\}$, $E =$ number of experts (Q -
 networks), $K =$ number of actions, $\Delta, \epsilon, \eta > 0$ and temperature $z, \tau > 0$, $n_r = -\infty$ (an upper
 bound on reward);

while True **do**

 Initialize episode by setting s_0

for $i = 1, 2, \dots, T$ (length of episode) **do**

 Observe state s_i ;

 Let probability of Q_k -network be $\rho_k = (1 - \eta) \frac{w_k}{\sum_{j=1}^E w_j} + \frac{\eta}{E}$;

 Sample network \bar{k} according to $\{\rho_k\}_k$;

 For $Q_{\bar{k}}$ -network, use ϵ -greedy to sample an action: $a^* = \arg \max_a Q_{\bar{k}}(s_i, a), j \in$
 $\{1, 2, \dots, K\}, \pi_j = (1 - \epsilon) \cdot \mathbb{1}_{j=a^*} + \frac{\epsilon}{K-1} \cdot \mathbb{1}_{j \neq a^*}$;

 Sample action a_i based on π ;

 Interact with the environment to receive reward r_i and next state s_{i+1} ;

$n_r = \max\{r_i, n_r\}$;

 Update the trust coefficient w_k of each Q_k -network as follows: $P_k = \epsilon$ -greedy(Q_k), $\hat{x}_{kj} =$

$1 - \frac{\mathbb{1}_{j=a^*}}{P_{kj} + \Delta} (1 - \frac{r_i}{n_r}), \forall j, y_k = E[\hat{x}_{kj}], w_k = w_k \cdot e^{\frac{y_k}{z}}$;

 Store (s_i, a_i, r_i, s_{i+1}) in experience replay buffer B ;

end for

 Update each expert's Q_k -network from buffer B

end while

286 idea is the same as in EXP4.P by using the experts that give advice vectors with deep Q -networks. It
 287 is a combination of deep neural networks with EXP4.P updates. From a different point of view, we
 288 can also view it as an ensemble in classification [20], by treating Q -networks as ensembles in RL.
 289 While general experts can be used, these are natural in a DQN framework. In our implementation
 290 and experiments we use two experts, thus $E = 2$ with two Q -networks. The first one is based on
 291 RND [5] while the second one is a simple DQN. To this end, in the algorithm before storing to the
 292 buffer, we also record $c_r^i = \|\hat{f}(s_i) - f(s_i)\|^2$, the RND intrinsic reward as in [5]. This value is
 293 then added to the 4-tuple pushed to B . When updating Q_1 corresponding to RND at the end of an
 294

295 iteration in the algorithm, by using $r_j + c_r^j$ we modify the Q_1 -network and by using c_r^j an update
 296 to \hat{f} is executed. Network Q_2 pertaining to ϵ -greedy is updated directly by using r_j . Intuitively,
 297 Algorithm 3 circumvents RND’s drawback with the total exploration guided by two experts with
 298 EXP4.P updated trust coefficients. When the RND expert drives high exploration, its trust coefficient
 299 leads to a high total exploration. When it has low exploration, the second expert DQN should have
 300 a high one and it incentivizes the total exploration accordingly. Trust coefficients are updated by
 301 reward estimates iteratively as in EXP4.P, so they keep track of the long-term performance of experts
 302 and then guide the total exploration globally. These dynamics of EXP4.P combined with intrinsic
 303 rewards guarantee global exploration. The experimental results exhibited in the next section verify
 304 this intuition regarding exploration behind Algorithm 3.

305 We point out that potentially more general RL algorithms based on Q -factors can be used, e.g., boost-
 306 rapped DQN [13], random prioritized DQN [12] or adaptive ϵ -greedy VDBE [19] are a possibility.
 307 Furthermore, experts in EXP4 can even be policy networks trained by PPO [15] instead of DQN for
 308 exploration. A recommendation is to have a good enough expert and a small number of experts.

309 5 Computational Study

310 As a numerical demonstration of the superior performance and exploration incentive of Algorithm 3,
 311 we show the improvements on baselines on two hard-to-explore RL games, Mountain Car and
 312 Montezuma’s Revenge. More precisely, we present that the real reward on Mountain Car improves
 313 significantly by Algorithm 3 in Section 5.1. Then we implement Algorithm 3 on Montezuma’s
 314 Revenge and show the growing and remarkable improvement of exploration in Section 5.2. Intrinsic
 315 reward $c_r^i = \|\hat{f}(s_i) - f(s_i)\|^2$ given by intrinsic model \hat{f} represents the exploration of RND in [5]
 316 as introduced in Sections 2 and 4. We use the same criterion for evaluating exploration performance
 317 of our algorithm and RND herein. RND incentivizes local exploration with the single step intrinsic
 318 reward but with the absence of global exploration.

319 5.1 Mountain Car

320 In this part, we summarize the experimental results of Algorithm 3 on Mountain Car, a classical
 321 control RL game. This game has very sparse positive rewards, which brings the necessity and
 322 hardness of exploration. Blog post [14] shows that RND based on DQN improves the performance of
 323 traditional DQN, since RND has intrinsic reward to incentivize exploration. We use RND on DQN
 324 from [14] as the baseline and show the real reward improvement of Algorithm 3, which supports the
 325 intuition and superiority of the algorithm.

326 The comparison between Algorithm 3 and RND is presented in Figure 1. Here the x-axis is the
 327 epoch number and the y-axis is the cumulative reward of that epoch. Figure 1a shows the raw
 328 data comparison between EXP4-RL and RND. We observe that though at first RND has several
 329 spikes exceeding those of EXP4-RL, EXP4-RL has much higher rewards than RND after 300 epochs.
 330 Overall, the relative difference of areas under the curve (AUC) is 4.9% for EXP4-RL over RND,
 331 which indicates the significant improvement of our algorithm. This improvement is better illustrated
 332 in Figure 1b with the smoothed reward values. Here there is a notable difference between EXP4-RL
 333 and RND. Note that the maximum reward hit by EXP4-RL is -86 and the one by RND is -118 ,
 334 which additionally demonstrates our improvement on RND.

335 We conclude that Algorithm 3 performs better than the RND baseline and that the improvement
 336 increases at the later training stage. Exploration brought by Algorithm 3 gains real reward on this
 337 hard-to-explore Mountain Car, compared to the RND counterpart (without the DQN expert). The
 338 power of our algorithm can be enhanced by adopting more complex experts, not limited to only DQN.

339 5.2 Montezuma’s Revenge and Pure exploration setting

340 In this section, we show the experimental details of Algorithm 3 on Montezuma’s Revenge, another
 341 notoriously hard-to-explore RL game. The benchmark on Montezuma’s Revenge is RND based on
 342 DQN which achieves a reward of zero in our environment (the PPO algorithm reported in [5] has
 343 reward 8,000 with many more computing resources; we ran the PPO-based RND with 10 parallel
 344 environments and 800 epochs to observe that the reward is also 0), which indicates that DQN has
 345 room for improvement regarding exploration.

346 To this end, we first implement the DQN-version RND (called simply RND hereafter) on Montezuma’s
 347 Revenge as our benchmark by replacing the PPO with DQN. Then we implement Algorithm 3 with

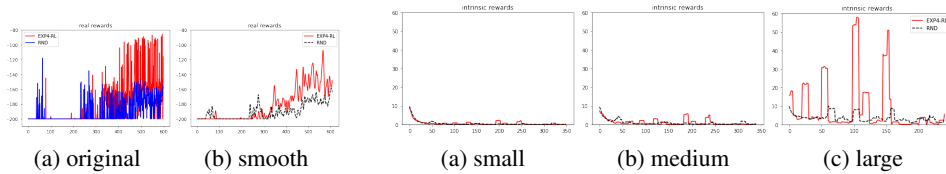


Figure 1: The performance of Algorithm 3 and RND measured by the epoch-wise re-intrinsic reward without parallel environments with three different burn-in periods on Mountain Car

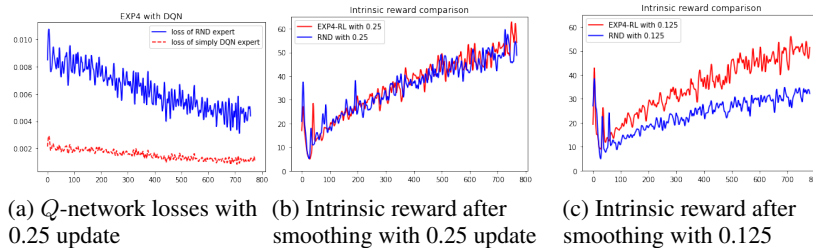


Figure 3: The performance of Algorithm 3 and RND with 10 parallel environments and with RND update probability 0.25 and 0.125, measured by loss and intrinsic reward.

two experts as aforementioned. Our computing environment allows at most 10 parallel environments. In subsequent figures the x-axis always corresponds to the number of epochs. RND update probability is the proportion of experience that are used for training the intrinsic model \hat{f} [5].

A comparison between Algorithm 3 (EXP4-RL) and RND without parallel environments (the update probability is 100% since it is a single environment) is shown in Figure 2 with the emphasis on exploration by means of the intrinsic reward. We use 3 different numbers of burn-in periods (58, 68, 167 burn-in epochs) to remove the initial training steps, which is common in Gibbs sampling. Overall EXP4-RL outperforms RND with many significant spikes in the intrinsic rewards. The larger the number of burn-in periods is, the more significant is the dominance of EXP4-RL over RND. EXP4-RL has much higher exploration than RND at some epochs and stays close to RND at other epochs. At some epochs, EXP4-RL even has 6 times higher exploration. The relative difference in the areas under the curves are 6.9%, 17.0%, 146.0%, respectively, which quantifies the much better performance of EXP4-RL.

We next compare EXP4-RL and RND with 10 parallel environments and different RND update probabilities in Figure 3. The experiences are generated by the 10 parallel environments.

Figure 3a shows that both experts in EXP4-RL are learning with decreasing losses of their Q -networks. The drop is steeper for the RND expert but it starts with a higher loss. With RND update probability 0.25 in Figure 3b we observe that EXP4-RL and RND are very close when RND exhibits high exploration. When RND is at its local minima, EXP4-RL outperforms it. Usually these local minima are driven by sticking to local maxima and then training the model intensively at local maxima, typical of the RND local exploration behavior. EXP4-RL improves on RND as training progresses, e.g. the improvement after 550 epochs is higher than the one between epochs 250 and 550. In terms for AUC, this is expressed by 1.6% and 3.5%, respectively. Overall, EXP4-RL improves RND local minima of exploration, keeps high exploration of RND and induces a smoother global exploration.

With the update probability of 0.125 in Figure 3c, EXP4-RL almost always outperforms RND with a notable difference. The improvement also increases with epochs and is dramatically larger at RND's local minima. These local minima appear more frequently in training of RND, so our improvement is more significant as well as crucial. The relative AUC improvement is 49.4%. The excellent performance in Figure 3c additionally shows that EXP4-RL improves RND with global exploration by improving local minima of RND or not staying at local maxima.

Overall, with either 0.25 or 0.125, EXP4-RL incentivizes global exploration on RND by not getting stuck in local exploration maxima and outperforms RND exploration aggressively. With 0.125 the improvement with respect to RND is more significant and steady. This experimental evidence verifies our intuition behind EXP4-RL and provides excellent support for it. With experts being more advanced RL exploration algorithms, e.g. DORA, EXP4-RL can bring additional possibilities.

383 **References**

- 384 [1] A. Agarwal, H. Luo, B. Neyshabur, and R. E. Schapire. Corraling a band of bandit algorithms.
385 In *Conference on Learning Theory*, pages 12–38. PMLR, 2017.
- 386 [2] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem.
387 *Machine learning*, 47(2-3):235–256, 2002.
- 388 [3] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit
389 problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.
- 390 [4] M. Bellemare, S. Srinivasan, G. Ostrovski, T. Schaul, D. Saxton, and R. Munos. Unifying
391 count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing*
392 *Systems*, pages 1471–1479, 2016.
- 393 [5] Y. Burda, H. Edwards, A. Storkey, and O. Klimov. Exploration by random network distillation.
394 In *International Conference on Learning Representations*, 2018.
- 395 [6] L. Fox, L. Choshen, and Y. Loewenstein. Dora the explorer: directed outreaching reinforcement
396 action-selection. In *International Conference on Learning Representations*, 2018.
- 397 [7] S. Grünewälder, J. Y. Audibert, M. Opper, and J. Shawe-Taylor. Regret bounds for gaussian
398 process bandit problems. In *Proceedings of the Thirteenth International Conference on Artificial*
399 *Intelligence and Statistics*, pages 273–280, 2010.
- 400 [8] X. Liu, F. Xia, R. L. Stevens, and Y. Chen. Cost-effective online contextual model selection.
401 *arXiv preprint arXiv:2207.06030*, 2022.
- 402 [9] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller.
403 Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- 404 [10] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves,
405 M. Riedmiller, A. K. Fidjeland, G. Ostrovski, and S. Petersen. Human-level control through
406 deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- 407 [11] G. Neu. Explore no more: Improved high-probability regret bounds for non-stochastic bandits.
408 *Advances in Neural Information Processing Systems*, 28, 2015.
- 409 [12] I. Osband, J. Aslanides, and A. Cassirer. Randomized prior functions for deep reinforcement
410 learning. In *Advances in Neural Information Processing Systems*, pages 8617–8629, 2018.
- 411 [13] I. Osband, C. Blundell, A. Pritzel, and B. Van Roy. Deep exploration via bootstrapped dqn. In
412 *Advances in Neural Information Processing Systems*, pages 4026–4034, 2016.
- 413 [14] O. Rivlin. Mountaincar_dqn_rnd. [https://github.com/orrivlin/MountainCar_DQN_](https://github.com/orrivlin/MountainCar_DQN_RND)
414 [RND](https://github.com/orrivlin/MountainCar_DQN_RND), 2019.
- 415 [15] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
416 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 417 [16] N. Srinivas, A. Krause, S. Kakade, and M. Seeger. Gaussian process optimization in the bandit
418 setting: no regret and experimental design. In *Proceedings of the 27th International Conference*
419 *on Machine Learning*, 2010.
- 420 [17] B. C. Stadie, S. Levine, and P. Abbeel. Incentivizing exploration in reinforcement learning with
421 deep predictive models. *arXiv preprint arXiv:1507.00814*, 2015.
- 422 [18] A. L. Strehl and M. L. Littman. An analysis of model-based interval estimation for markov
423 decision processes. *Journal of Computer and System Sciences*, 74(8):1309–1331, 2008.
- 424 [19] M. Tokic. Adaptive ϵ -greedy exploration in reinforcement learning based on value differences.
425 In *Annual Conference on Artificial Intelligence*, pages 203–210. Springer, 2010.
- 426 [20] R. Xia, C. Zong, and S. Li. Ensemble of feature sets and classification algorithms for sentiment
427 classification. *Information Sciences*, 181(6):1138–1152, 2011.
- 428 [21] S. Zuo. Near optimal adversarial attack on UCB bandits. *arXiv preprint arXiv:2008.09312*,
429 2020.