

1 We would like to thank the reviewers for their time. We humbly request for more of your time to reevaluate scores,  
2 given our responses below. We propose a general procedure to modify bandit algorithms for stochastic contextual bandit  
3 problems to be compatible with multiple adversarial bandit algorithms which allows us to obtain previously unattainable  
4 model selection regret guarantees that can be applied to a wide variety of problems in the setting of i.i.d. contexts. These  
5 assumptions (stochastic environments and i.i.d. contexts) are fairly standard in the literature. Because of its portability,  
6 this allows us to easily plug in multiple existing bandit algorithms and obtain regret guarantees for many problems for  
7 which no model selection guarantees were known before. Notably, we are able to obtain meaningful model selection  
8 guarantees even when the best base algorithm's regret is not fully known. We summarize our contributions:

9 **(Section 4.1) Mis-specified contextual bandit with unknown error:** We provide the first solution for the case of  
10 changing action sets. For fixed action sets, we improve upon the best existing result (Lattimore et al, 2020), and match  
11 the lower bound. To understand this result, consider  $k$  arms of dimension  $d < \sqrt{k}$ . If the arms are linear, then UCB has  
12 regret  $\tilde{O}(\sqrt{kT})$  and LinUCB has  $\tilde{O}(d\sqrt{T})$ , so the best base's regret is  $\tilde{O}(d\sqrt{T})$ . If the arms are not linear, then UCB  
13 has regret  $\tilde{O}(\sqrt{kT})$  and LinUCB has linear regret, so the best base's regret is  $\tilde{O}(\sqrt{kT})$ . It is impossible to know the  
14 best base's regret without knowing the environment, but our method achieves a regret matching the lower bound.

15 **(Section 4.2) Linear contextual bandit with unknown dimension:** For finite action sets, we provide a regret bound  
16 that does not depend on the number of actions, in contrast to the best existing result (Foster et al, 2019). We provide the  
17 first solution for infinite action sets. We provide the first solution when both the dimension and the mis-specified error  
18 are unknown. This is also the first result in literature that can combine multiple types of model selection.

19 **(Section 4.2) Non-parametric contextual bandit with unknown dimension:** we provide the first solution.

20 **(Section 4.3) Tuning the exploration rate of  $\epsilon$ -greedy:** we provide the first solution for this problem.

21 **(Section 4.4) Choosing between multiple feature maps for RL:** we provide the first solution.

22 **(Appendix A1) Generalized linear bandits with unknown link function:** we provide this problem's first solution.

23 **(Appendix A2) Bandit with heavy tail:** we provide the first solution for this problem.

24 **Reviewer 1:** "The selection of the range": The regret is multiplied by at most a factor of the number of bases  $M$ , which  
25 is  $\log(B)$  if the upper bound of the parameter is  $B$ . Therefore  $B$  can be chosen quite loosely as long as  $\log(B)$  is not  
26 too large. In the paper we choose the largest  $\epsilon$  to be 100,000 but in practice such a large  $\epsilon$  is unreasonable. Still, even  
27 with such a loosely chosen  $B$ , the performance is reasonably well.

28 We would like to emphasize that the exponential grid division of the parameter space is not central to our key  
29 contributions. How to define the best range and search efficiently in the parameter space is an interesting but separate  
30 research topic not in the scope of this paper.

31 **Reviewer 2:** Thank you for your comments. The citation to the non-stationary bandit paper was left from an earlier  
32 version of this submission, and will be removed.

33 **Reviewer 3: Weaknesses:** A) Relation to prior work in Section 4: In Section 4, we show that the strategy can be  
34 seamlessly used to solve a number of challenging open problems (reiterated above). We provide the first solution for  
35 these problems, and there are no existing solutions to compete with.

36 B) We would present the full description of the algorithm before Section 4.

37 C) "The assumption that Base algorithms only have access to rewards of rounds when they are selected...": This is not  
38 an assumption but a standard algorithmic design choice (modularity). It is not restrictive because we are able to produce  
39 a variety of state-of-the-art results (listed at the beginning). It is not clear if letting the base algorithms share the rewards  
40 will increase the performance.

41 **Correctness:** A) "I think it's not trivial to reproduce the results...": We would like more explanations as to why it would  
42 be difficult to reproduce the results. We believe we have provided all the details for our experiments. We have listed all  
43 configurations, parameter choices and number of repetitions. We also reproduced the master algorithms in Appendix B  
44 for convenience.

45 B) "...the instantaneous regret of Step 2 is  $1/s$  times...": Let the cumulative regret of step 1 at round  $s$  be  $U(s)$ , then the  
46 cumulative regret of step 2 at round  $S$  is:  $\frac{1}{1}U(1) + \frac{1}{2}U(2) + \dots + \frac{1}{S}U(S) < (\frac{1}{1} + \frac{1}{2} + \dots + \frac{1}{S})U(S) = \log(S)U(S)$ .  
47 More details are in the proof in Lemma F1, page 24.

48 **Additional feedback:** A) "On line 94: what is  $i$ ?":  $i$  is the index of the base.

49 B) "...cases where the action set changes...": In Step 2 of the smoothing procedure, we repeat the policy (which is the  
50 way to compute the chosen action), not the action. Even if the action set changes, the policy stays the same.

51 C) "...lines 6 through 11 in the pseudocode...": this is for the analysis of term II, as explained from line 90-98.