

1 We thank all reviewers for their detailed feedback.

2 1 Reviewer 1

- 3 • About the solution for under-exploration: We agree that uniform exploration may not be the best possible
4 approach. However, we think it makes the important point that adding some exploration in response to
5 approximate inference can guarantee sub-linear regret, even though this particular form of exploration might
6 not be optimal.
- 7 • About stochastic optimism: Thank you for the reference, we agree that stochastic optimism in RL might be
8 related. At first glance, stochastic optimism looks quite different from a bound on the error of alpha-divergence,
9 but it would definitely be interesting to see if there is a deeper connection.
- 10 • About experiments: We agree that Q_t and Z_t (which serve to illustrate the intuition about over-exploration
11 and under-exploration) might not be natural approximations. However Figure 3 shows variational inference
12 and ensemble sampling on a 50 armed bandit instance, which is more realistic.

13 2 Reviewer 2

- 14 • About Garbage-in Reward-out paper [Kveton, 2019]: Thank you for the reference. This paper is related but
15 does not study the same problem – they are not concerned with the case where the Thompson sampling’s
16 inference oracle is approximate. We would cite and discuss it more thoroughly in a revision.
- 17 • About presentation issues: Thank you very much for the comments. We will revise accordingly.
- 18 • About the implication of Theorem 1: Typically, approximate inference methods minimize divergences. Broadly
19 speaking, we show that making a divergence a small constant, alone, is not enough to guarantee sub-linear
20 regret. We do not mean to imply that low regret is *impossible* but simply that making an alpha-divergence a
21 small constant alone is not sufficient. We will clarify this point.
- 22 • About the implications of all theorem statements:
23 Broadly speaking, Theorem 1 implies that when the approximation scheme over-explores, even though the
24 posterior may concentrate, we suffer regret because the approximation chooses the sub-optimal arm with
25 higher probability than the posterior at every time-step due to over-exploration.
26 On the other hand, Theorem 2 implies that when the approximation scheme under-explores, the posterior may
27 not concentrate and therefore chooses the sub-optimal arm most of the times, leading the approximation to do
28 the same. Theorem 3 strengthens Theorem 2’s observation by showing that adding forced exploration will
29 help the posterior to concentrate and choose the optimal arm most of the times, leading the under-explored
30 approximation to do the same.
- 31 • About the correctness of the theorem statements: We believe the results are correct. It would be great if the
32 reviewer could point out the points of concern.

33 3 Reviewer 3

- 34 • We apologize for the grammatical mistakes and thank you so much for listing them.
- 35 • Thank you also for the suggested edits about Bayesian agent in line 52 and reverse KL divergence in line 76.
36 We will revise accordingly.