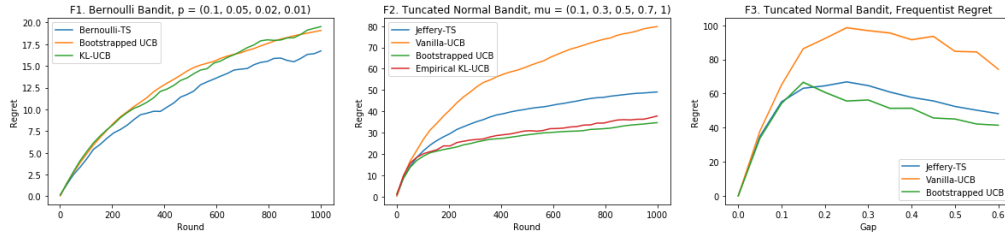1 We would like to thank all reviewers for your valuable and detailed comments! We think we can fix all the issues raised
2 by the reviewers and add more experimental results in the final version of the paper. We hope you are satisfied with our
3 point-by-point responses and increase your scores.

4 **Response to Reviewer 1**

5 *Experiments.* Thank you for pointing out KL-UCB, a great and powerful algorithm for cases with bounded rewards. We
6 will compare with KL-UCB in detail, and have done comprehensive experiments with KL-UCB, using Prof. Olivier
7 Cappé's package online (due to the space limit, we can only report parts of them here). Our algorithm is comparable
8 with KL-UCB (Bernoulli, F1) and empirical KL-UCB (truncated normal, F2) under fixed problem instances suggested
9 in the package. Although KL-UCB is also data-dependent, our proposed method is from a very different non-parametric
10 perspective and uses different tools by bootstrap. In practice, resampling tends to be more efficient computationally,
11 without solving a convex optimization each round like empirical KL-UCB. Moreover, our method can work with
unbounded rewards and we believe it is easier to generalize to structured bandits, e.g. linear bandit.

12



13 The UCB1 we use, defined as $\bar{y}_{n_{k,t}} + \sigma\sqrt{2\log(t)/n_{k,t}}$, exactly follows your suggestion (see E.1 in the appendix). So
14 we will change its name to vanilla UCB. In addition, we have added the frequentist regret curve (F3). Here, the regrets
15 of various algorithms are with respect to the instance gap ($\Delta$) and $\mu = (\Delta, 0, \ldots, 0)$. In the linear bandit part, the
16 dimension is specified in the title of each figure.

17 *Complexity.* Indeed, our algorithm is more complex than vanilla UCB and requires more memories. The computational
18 complexity at step $t$ is $\widetilde{\mathcal{O}}(Bt) \le \widetilde{\mathcal{O}}(BT)$. Comparing with vanilla UCB, the extra $Bt$ is due to resampling. We also
19 derive Theorem 0.1 below for MC quantile approximation error that provides us a theoretical guidance for the selection
20 of $B$. In practice, the choice of $B$ is seldom treated as a tuning parameter, but usually determined by the available
21 computational resource. To reduce the computational cost, we could use the idea of Bag of Little Bootstraps [2]. For
22 $\delta$, from (2.5), a smaller value of $\delta$ enables us to calculate the quantile at a closer level to $\alpha$ but will result in a larger
23 correction term. Since the correction term converges to 0 faster, we suggest a smaller $\delta$. In Section 4.2, the $\delta = 1/(1+t)$
24 is essentially the confidence level, rather than a hyper-parameter $\delta$ in (2.5). The choice of $1/T^2$ led to an easy analysis.
25 Using similar techniques in Chapter 8.2 of [1], we can derive a similar regret bound by setting $\alpha_t = 1/(t\log^\alpha(t))$ for
26 any $\alpha > 0$. We re-run the experiments and there is no significant change.

27 *Proof Clarifications.* Thank you for pointing out the typos and notations! We have revised them accordingly. For
28 Equation B.2, by the symmetric assumption of the reward, the distribution of $y_i - \mu$ is *exactly the same* as the distribution
29 of $w_i(y_i-\mu)$ for Rademacher r.v. $\{w_i\}$. This implies $\mathbb{P}_{\boldsymbol{y}}(\frac{1}{n}\sum_{i=1}^n(y_i-\mu) > q_\alpha(\boldsymbol{y}_n-\mu)) = \mathbb{P}_{\boldsymbol{y},\boldsymbol{w}}(\frac{1}{n}\sum_{i=1}^n w_i(y_i-\mu) >$
30 $q_\alpha((\boldsymbol{y}_n - \mu) \circ \boldsymbol{w}_n)) = \mathbb{E}_{\boldsymbol{w}}\mathbb{P}_{\boldsymbol{y}}\left(\frac{1}{n}\sum_{i=1}^n w_i(y_i - \mu) > q_\alpha((\boldsymbol{y}_n - \mu) \circ \boldsymbol{w}_n)\right)$. For the second one, instead of using
31 conditional event, we could use union event trick: $\mathbb{P}(A) = \mathbb{P}(A \cap B) + \mathbb{P}(A \cap B^c) \le \mathbb{P}(A \cap B) + \mathbb{P}(B^c)$. By choosing
32 $A = \{\bar{y}_n - \mu > q_{\alpha(1-\delta)}(\boldsymbol{y}_n - \bar{y}_n) + (2\log(2/\alpha\delta)/n)^{1/2}\varphi(\boldsymbol{y}_n)\}$ and $B = \{\boldsymbol{y}_n \in \mathcal{E}\}$, we can reach the conclusion.

33 **Response to Reviewer 2**

34 (1)(2). We will explicitly clarify and mention them in the beginning. (3). Yes, you are right. (4). We can derive a similar
35 regret bound by setting $\alpha_t = 1/(t\log^\alpha(t))$ for any $\alpha > 0$. Please see our response for complexity part above. (5). We
36 will change the argument to heavier tail than sub-Gaussian/exponential. (6). Thanks for the detailed references. We will
37 add them in the introduction.

38 **Response to Reviewer 3**

39 Thank you for pointing the MC quantile approximation. We have derived the corresponding theorem for the control of
40 the approximation of the bootstrapped quantile.

41 **Theorem 0.1** (Monte Carlo Quantile Approximation). Suppose the same conditions in Theorem 2.2 hold. We have
42 $\mathbb{P}_{\boldsymbol{y},\boldsymbol{w}}(\bar{y}_n - \mu > \widetilde{q}_\alpha^B + \sqrt{\log(2/\alpha\delta)/n}\varphi(\boldsymbol{y}_n)) \le \alpha + \frac{\lfloor B\alpha \rfloor + 1}{B+1} \le 2\alpha + \frac{1}{B+1}$, where $\widetilde{q}_\alpha^B$ is the Monte Carlo approximated
43 quantile defined in (D.1).

44 By replacing the true quantile $q_\alpha$ by a MC quantile $\widetilde{q}_\alpha^B$ based on $B$ i.i.d bootstrapped weights, we lose at most $1/(B+1)$
45 for the confidence level. The proof is similar to Theorem 2.2 except for the control of i.i.d approximation error.

46 [1]. Bandit Algorithms. Cambridge University Press (2019). [2]. The Big Data Bootstrap. ICML (2012).