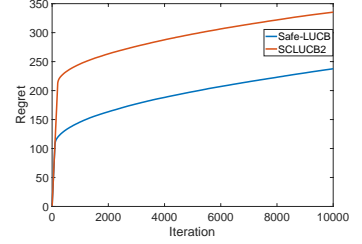


1 We thank the reviewers for their constructive feedback and their valuable time.

2 **Reviewer 1:** First, we respond to the reviewer’s **questions/suggestions regard-**
 3 **ing the experimental results.** Regarding comparisons to the safe-LUCB of [16],
 4 we present SCLUCB2 in App. H as a modified version of our main SCLUCB
 5 algorithm tailored for the exact safe bandit setting studied in [16]. Importantly,
 6 SCLUCB2 comes with a theoretical regret bound, which matches the proposed
 7 problem-dependent upper bound in [16]. We now confirm this numerically in the
 8 displayed figure, which plots the cumulative regret of the two algorithms averaged
 9 over 100 realizations. We will include this new numerical study in the final version.



10 The reason that we only plot regret curves and not the number of times the safety constraint is violated, is because this
 11 number is *zero* for almost all realizations. This is expected since all our algorithms guarantee the model’s requirement
 12 that the safety constraints are not violated for any time step, with high-probability $1 - \delta$. Second, **regarding the**
 13 **parameters (R, S, L) of Ass. 1-3**, assuming knowledge of them is standard in the literature of linear bandits (see
 14 [5,11,10,13-15]). Their specific values are, of course, highly application-dependent, but the underlying hypothesis is
 15 that they can be accurately determined based on domain-knowledge/physics, or, estimated from historic data. Even
 16 if accurate approximations are not possible, rather loose bounds suffice to run the algorithms. Of course, the quality
 17 of these bounds affects the performance, but, the accompanying regret-bounds quantify the effect. **Regarding the**
 18 **parameters $(r_l, r_h, \kappa_l, \kappa_h)$ in Ass. 4** that are associated with the baseline policy, it can be reasonably assumed that they
 19 can be estimated accurately from data. This is because we think of the baseline policy as “past strategy”, implemented
 20 before bandit-optimization, thus producing large amount of data (see also [1-3]). If no knowledge is available however,
 21 κ_h and r_h can always be set to equal 1 (since for simplicity we assume that the mean rewards are in $[0, 1]$). Similarly, κ_l
 22 can be set equal to zero. On a related note, we address the **question on tuning the hyper-parameters $\delta, \lambda, \rho, \alpha$** .
 23 The tuning of δ, λ is standard and is same as in all linear-bandit algorithms [5,11,10,13-15]: $1 - \delta \in (0, 1)$ is the desired
 24 confidence (e.g., 0.95) on the algorithm’s realizations to satisfy the regret bounds (here, also the safety constraints); the
 25 regularization parameter λ can be set equal to one. The parameter ρ , controlling the exploration level of conservative
 26 actions can take any value in the interval specified in Lemma 2.2. The parameter $\alpha \in (0, 1)$, controlling the conservatism
 27 level of the learning process, is assumed known to the learner similar to [1,2,3]. We will clarify the above in the
 28 revision. Finally, **the assumption $\langle x, \theta_* \rangle \leq 1$** is not essential and is rather only meant for simplicity. Specifically, the
 29 assumptions $\|x\|_2 \leq L$ and $\|\theta_*\|_2 \leq S$ suffice, as they guarantee the constant bound LS for $\langle x, \theta_* \rangle$; thus, nothing
 30 fundamental changes in our analysis. For example, without this assumption, ρ_3 in Eq. (18) of Theorem 4.1. simply
 31 changes to $\rho_3 = \frac{\alpha r_l}{S+LS}$. Contrary to our intention, this assumption appears to be confusing and will be removed in the
 32 final version. Minor: The parameter ρ in Algo. 1 appears in the definition of x_t^{cb} in Eq. (11). We will clarify this.

33 **Reviewer 2:** First, please refer to lines 18-22 above on how the parameters $(r_l, r_h, \kappa_l, \kappa_h)$ are chosen. We further
 34 clarify the following. Regarding κ_l : Indeed, there is a typo in line 228 and the related factor in the sample complexity
 35 should rather be $\kappa_l + \alpha r_l$ as specified explicitly in Thm. 3.3. What this bound suggests is that while setting $\kappa_l = 0$
 36 is possible, a higher value is preferable (provided that it lower bounds κ_{b_t} in (4)), since it results in smaller regret.
 37 Regarding the requirement $r_l > 0$: Indeed, this is necessary for the algorithms to perform well and is critically
 38 used in the proofs (e.g. Eq. (23)). That said, this is expected to be met in practice since the baseline policy is the
 39 system’s current strategy and should have been associated with at least a positive reward. Second, let us clarify the
 40 **assumption that the action set contains the unit ball**, eqv. $L \geq 1$. This goes hand-in-hand with our assumption
 41 $\|\zeta\|_2 = 1$ (line 199), since together they guarantee that the convex combination x_t^{cb} in (11) is feasible, i.e., satisfies
 42 $\|x_t^{\text{cb}}\|_2 \leq L$. However, this requirement remains true as long as $\|\zeta\|_2 = \epsilon$ and $L \geq \epsilon$ for any $\epsilon > 0$. To see this, note
 43 that $\|x_t^{\text{cb}}\|_2 \leq (1 - \rho)L + \rho\epsilon \leq L - \rho(L - \epsilon) \leq L$. In particular, $\epsilon > 0$ can be chosen small enough for “thin sets”.
 44 Changing $\|\zeta\|_2 = \epsilon$, we simply adjust $h_1 = 2\rho_1(1 - \rho_1)L\epsilon + 2\rho_1^2\epsilon^2$ and $\rho_1 = \frac{\alpha r_l}{S\epsilon + r_h}$ in Thm. 3.3. Minor: Thank you
 45 for the comment about LUCB. Also, we agree and will modify line 59 to clarify that the learner Knows x_{b_t} .

46 **Reviewer 3:** Thank you for the suggestion on **numerically verifying the number of**
 47 **times that baseline is played.** The figure on the right plots the cumulative number
 48 of baseline actions played by SCLTS until time t , for $t = 1, \dots, 1000$. The solid line
 49 depicts average over 100 realizations and the shaded regions show standard deviation.
 50 The figure confirms the logarithmic trend predicted by theory. We will upload our code
 51 as suggested. We finish with **a brief proof-sketch of Thm. 3.3**, which we will include
 52 in the paper. The first idea is based on the intuition that if a baseline action is played at
 53 around t , then the algorithm does not yet have a good estimate of the unknown parameter
 54 θ_* and the safe actions played thus far have not yet expanded properly in all directions.
 55 Formally, this translates to small $\lambda_{\min}(V_t)$ and the *upper* bound $O(\log \tau) \geq \lambda_{\min}(V_\tau)$ (Eq. (43)). The second key idea
 56 is to exploit the randomized nature of the conservative actions (cf. (11)) to *lower* bound $\lambda_{\min}(V_\tau)$ by the number $|N_\tau^c|$
 57 of times that SCLTS plays the baseline actions up to that round (cf. Lemma D.1). Putting these together leads to the
 58 advertised upper bound $O(\log T)$ on the total number $|N_T^c|$ of times the algorithm plays the baseline actions.

