

1 Dear reviewers, first of all we would like to thank you for taking the time to review our paper during those challenging  
2 times! Answers to your questions are in place.

3 **Exponential family of stochastic matrices:** We sketch two examples illustrating that one-parameter exponential  
4 families of stochastic matrices generalize all the applications of exponential families of discrete probability distributions.

5 1. Take example 1 from the paper. As discussed in the paper this generalizes the Bernoulli exponential family. In  
6 the i.i.d. bandit model one would have  $K$  Bernoulli processes,  $\text{Ber}(p_{\theta_1}), \dots, \text{Ber}(p_{\theta_K})$ . In the Markovian bandit  
7 model under consideration one has  $K$  Markov processes specified by  $K$  stochastic matrices  $P_{\theta_1}, \dots, P_{\theta_K}$ ,  
8 and each row of them is specified by a coin flip as discussed in example 1. Those  $K$  Markovian processes  
9 form a generalization of the  $K$  Bernoulli processes.

10 2. In the same spirit as before, we will sketch an approximation for the generalization of the Poisson exponential  
11 family, which is useful for count data. Due to the finite state space, approximate the  $\text{Pois}(\lambda)$  distribution,  
12 with  $\text{Bin}(n, \lambda/n)$ , and  $n$  large enough so that the two distributions are  $\epsilon$ -close with respect to the total  
13 variation distance. Now our state space is  $S = \{0, \dots, n\}$ . For the generator stochastic matrix we pick  
14  $n + 1$  row distributions  $\text{Bin}(n, \lambda_0/n), \dots, \text{Bin}(n, \lambda_n/n)$  which are approximately  $\text{Pois}(\lambda_0), \dots, \text{Pois}(\lambda_n)$ ,  
15 and for the Markovian bandit model we tilt the generator stochastic matrix by parameters  $\theta_1, \dots, \theta_K$ , to obtain  
16  $K$  Markovian processes, each of them giving rewards and transitioning according to approximate Poisson  
17 distributions. In the i.i.d. case we would have just a single  $\text{Pois}(\lambda)$  distribution, and we would produce  $K$  tilts  
18 which would correspond to the distributions of the  $K$  arms.

19 The problem with generalizing even further to countably infinite state spaces, or even continuous state spaces is the  
20 peculiar behavior of eigenvalues and eigenfunctions (which may not even exist) on infinite dimensional spaces.

21 **Round-robin scheme:** As discussed in the paper the round-robin idea dates back to Lai and Robbins, although forgotten  
22 nowadays. In this paper:

23 1. We use different statistics to make the scheme computable in the case of multiple-plays and Markovian rewards  
24 (note that coming up with a computable algorithm in the case of multiple-plays is the motivation of [18]  
25 Komiyama, Honda, Nakagawa, where they study Thompson sampling for multiple plays and i.i.d. Bernoulli  
26 rewards).

27 2. We provide a finite-time analysis (as opposed to asymptotics in prior work).

28 3. Through experimental results we bring to the attention of the research community that this type of scheme  
29 is computationally more efficient, and equally as effective as the status quo of calculating UCB or KL-UCB  
30 indices for each arm at each round. In particular for KL-UCB type of indices the computational improvement  
31 can be quantified as  $O(K + \log 1/\epsilon)$  (this paper) vs  $O(K \log 1/\epsilon)$  (KL-UCB paper) cost per round, where  $K$   
32 is the number of arms, and  $\epsilon$  is the quality of the KL divergence inverses that we're interested in.

33 Reviewer #3 we don't think that your suggestion as an alternative to the round-robin scheme works. For instance take  
34  $B = T/2$ , where  $T$  is the time horizon. Then you're claiming that by calculating UCB scores just twice, the regret  
35 might only get double, which is clearly wrong. Additionally, the generator stochastic matrix  $P$  and the function  $f$  need  
36 not be known for the algorithm. All that needs to be known is what is declared as parameters in the preamble of the algorithm, so  
37 in particular the only thing that needs to be known from the family is the KL divergence rate function. Finally, we could even  
38 eliminate the presence of  $\delta$ , but we decided to keep it as a knob to tune the algorithm. In particular for the experiments  
39 we tuned  $\delta$  by playing around with several values in the range  $(0, 1/K)$ .

40 **Maximal inequality:** The workhorse behind a maximal inequality is typically some variant of Doob's martingale  
41 inequality. For our paper we have reviewed the literature on martingale methods to derive deviation inequalities for  
42 Markov chains or more general dependent processes [8, 15, 19, 24, 25, 26, 27, 30, 33, 34], and to the best of our  
43 knowledge none of them seems capable to deliver a maximal inequality, due to the fact that they don't use the exponential  
44 martingale (lemma 1), but instead they use Doob's or Dynkin's martingale. For your interest, A. Kontorovich is one of  
45 the authors of [15] and also an AC, so maybe you could consult him for an extra opinion?