

A KL-LUCB algorithm for Large-Scale Crowdsourcing

Authored by:

Robert Nowak
Ervin Tanczos
Bob Mankoff

Abstract

This paper focuses on best-arm identification in multi-armed bandits with bounded rewards. We develop an algorithm that is a fusion of lil-UCB and KL-LUCB, offering the best qualities of the two algorithms in one method. This is achieved by proving a novel anytime confidence bound for the mean of bounded distributions, which is the analogue of the LIL-type bounds recently developed for sub-Gaussian distributions. We corroborate our theoretical results with numerical experiments based on the New Yorker Cartoon Caption Contest.

1 Paper Body

This paper focuses on best-arm identification in multi-armed bandits with bounded rewards. We develop an algorithm that is a fusion of lil-UCB and KL-LUCB, offering the best qualities of the two algorithms in one method. This is achieved by proving a novel anytime confidence bound for the mean of bounded distributions, which is the analogue of the LIL-type bounds recently developed for sub-Gaussian distributions. We corroborate our theoretical results with numerical experiments based on the New Yorker Cartoon Caption Contest.

1

Multi-Armed Bandits for Large-Scale Crowdsourcing

This paper develops a new multi-armed bandit (MAB) for large-scale crowdsourcing, in the style of the KL-UCB [4, 9, 3]. Our work is strongly motivated by crowdsourcing contests, like the New Yorker Cartoon Caption contest [10]3 . The new approach targets the “best-arm identification problem” [1] in the fixed confidence setting and addresses two key limitations of existing theory and algorithms: (i) State of the art algorithms for best arm identification are based on sub-Gaussian confidence bounds [5] and fail to exploit the fact that rewards are usually bounded in crowdsourcing applications. (ii) Existing KL-UCB algorithms for best-arm identification do exploit bounded rewards [8] , but

have suboptimal performance guarantees in the fixed confidence setting, both in terms of dependence on problem-dependent hardness parameters (Chernoff information) and on the number of arms, which can be large in crowdsourcing applications. The new algorithm we propose and analyze is called lil-KLUCB, since it is inspired by the lil-UCB algorithm [5] and the KL-LUCB algorithm [8]. The lil-UCB algorithm is based on sub-Gaussian bounds and has a sample complexity for best-arm identification that scales as $X \sum_{i=1}^n \log \frac{1}{\delta} \log \frac{1}{\Delta_i}$, where $\delta \in (0, 1)$ is the desired confidence and $\Delta_i = \mu_1 - \mu_i$ is the gap between the means of the best arm (denoted as arm 1) and arm i . If the rewards are in $[0, 1]$, then the KL-LUCB algorithm has $\sum_{i=1}^n \log \frac{1}{\Delta_i}$.

where $\delta \in (0, 1)$ is the desired confidence and $\Delta_i = \mu_1 - \mu_i$ is the gap between the means of the best arm (denoted as arm 1) and arm i . If the rewards are in $[0, 1]$, then the KL-LUCB algorithm has $\sum_{i=1}^n \log \frac{1}{\Delta_i}$.

This work was partially supported by the NSF grant IIS-1447449 and the AFSOR grant FA9550-13-1-0138. For more details on the New Yorker Cartoon Caption Contest, see the Supplementary Materials.

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

a sample complexity scaling essentially like $X \sum_{i=1}^n \log \frac{1}{\delta} \log \frac{1}{\Delta_i}$, where n is the number of arms and $\Delta_i := \mu_1 - \mu_i$ is the Chernoff-

information between a $\text{Ber}(\mu_1)$ and a $\text{Ber}(\mu_i)$ random variable⁵. Ignoring the logarithmic factor, this bound is optimal for the case of Bernoulli rewards [7, 11]. Comparing these two bounds, we observe that KL-LUCB may offer benefits since $\Delta_i = \mu_1 - \mu_i \geq (\mu_1 - \mu_i)^2 / 2 = \Delta_i^2 / 2$, but lil-UCB has better logarithmic dependence on the Δ_i and no explicit dependence on the number of arms n . Our new algorithm lil-KLUCB offers the best of both worlds, providing a sample complexity that scales essentially like $X \sum_{i=1}^n \log \frac{1}{\delta} \log \frac{1}{\Delta_i}$.

The key to this result is a novel anytime confidence bound for sums of bounded random variables, which requires a significant departure from previous analyses of KL-based confidence bounds. The practical benefit of lil-KLUCB is illustrated in terms of the New Yorker Caption Contest problem [10]. The goal of that crowdsourcing task is to identify the funniest cartoon caption from a batch of $n \approx 5000$ captions submitted to the contest each week. The crowd provides 3-star ratings for the captions, which can be mapped to $\{0, 1/2, 1\}$, for example. Unfortunately, many of the captions are not funny, getting average ratings close to 0 (and consequently very small variances). This fact, however, is ideal for KL-based confidence intervals, which are significantly tighter than those based on sub-Gaussianity and the worst-case variance of $1/4$. Compared to existing methods, the lil-KLUCB algorithm better addresses the two key features in this sort of application: (1) a very large number of arms, and (2) bounded reward distributions which, in many cases, have very low variance. In certain instances, this can have a profound effect on sample complexity (e.g., $O(n^2)$ complexity for algorithms using sub-Gaussian bounds vs. $O(n \log n)$ for lil-KLUCB, as shown in Table 1). The paper is organized as follows. Section 2 defines the best-arm identification problem, gives the lil-KLUCB algorithm and states the main results. We also briefly review related literature, and compare

the performance of lil-KLUCB to that of previous algorithms. Section 3 provides the main technical contribution of the paper, a novel anytime confidence bound for sums of bounded random variables. Section 4 analyzes the performance of the lil-KLUCB algorithm. Section 5 provides experimental support for the lil-KLUCB algorithm using data from the New Yorker Caption Contest.

2

Problem Statement and Main Results

Consider a MAB problem with n arms. We use the shorthand notation $[n] := \{1, \dots, n\}$. For every $i \in [n]$ let $\{X_{i,j}\}_{j \in \mathbb{N}}$ denote the reward sequence of arm i , and suppose that $P(X_{i,j} \in [0, 1]) = 1$ for all $i \in [n]$, $j \in \mathbb{N}$. Furthermore, assume that all rewards are independent, and that $X_{i,j} \leq \mu_i$ for all $j \in \mathbb{N}$. Let the mean reward of arm i be denoted by μ_i and assume w.l.o.g. that $\mu_1 \leq \mu_2 \leq \dots \leq \mu_n$. We focus on the best-arm identification problem in the fixed-confidence setting. At every time $t \in \mathbb{N}$ we are allowed to select an arm to sample (based on past rewards) and observe the next element in its reward sequence. Based on the observed rewards, we wish to find the arm with the highest mean reward. In the fixed confidence setting, we prescribe a probability of error $\delta \in (0, 1)$ and our goal is to construct an algorithm that finds the best arm with probability at least $1 - \delta$. Among $1 - \delta$ accurate algorithms, one naturally favors those that require fewer samples. Hence proving upper bounds on the sample complexity of a candidate algorithm is of prime importance. The lil-KLUCB algorithm that we propose is a fusion of lil-UCB [5] and KL-LUCB [8], and its operation is essentially a special instance of LUCB++ [11]. At each time step t , let $T_i(t)$ denote the total number of samples drawn from arm i so far, and let $\hat{\mu}_i(t)$ denote corresponding empirical mean. The algorithm is based on lower and upper confidence bounds of the following general form: 4

A more precise characterization of the sample complexity is given in Section 2. The Chernoff-information between random variables $\text{Ber}(x)$ and $\text{Ber}(y)$ ($0 \leq x \leq y \leq 1$) is $D(x, y) = \inf_{z \in [x, y]} D(z, x) = D(z, y)$, where $D(z, x) = z \log \frac{z}{xz} + (1-z) \log \frac{1-z}{1-x}$ and z^* is the unique $z \in [x, y]$ such that $D(z, x) = D(z, y)$. 5

2

for each $i \in [n]$ and any $\delta \in (0, 1)$

$$c \log(\delta \log 2(2T_i(t))/\delta) \leq L_i(t, \delta) = \inf_{\mu_i \leq \mu \leq \hat{\mu}_i(t)} D(\mu, \hat{\mu}_i(t))$$

$$c \log(\delta \log 2(2T_i(t))/\delta) \leq U_i(t, \delta) = \sup_{\mu_i \leq \mu \leq \hat{\mu}_i(t)} D(\mu, \hat{\mu}_i(t))$$
, $m \leq T_i(t)$ where c and δ are small constants (defined in the next section). These bounds are designed so that with probability at least $1 - \delta$, $L_i(T_i(t), \delta) \leq \mu_i \leq U_i(T_i(t), \delta)$ holds for all $t \in \mathbb{N}$. For any $t \in \mathbb{N}$ let $\text{TOP}(t)$ be the index of the arm with the highest empirical mean, breaking ties at random. With this notation, we state the lil-KLUCB algorithm and our main theoretical result. lil-KLUCB 1. Initialize by sampling every arm once. 2. While $L_{\text{TOP}(t)}(T_{\text{TOP}(t)}(t), \delta/(n-1)) >$

$\max_{i \neq \text{TOP}(t)} U_i(T_i(t), \delta)$ do:

$i_6 = \text{TOP}(t)$

δ Sample the following two arms: $i_6 = \text{TOP}(t)$, and $i_7 = \arg \max_{i \neq \text{TOP}(t)} U_i(T_i(t), \delta)$

and update means and confidence bounds. 3. Output $\text{TOP}(t)$ Theorem 1. For every $i \geq 2$ let $\mu_i = (\mu_i, \sigma_i)$, and $\mu = \max_{i \geq 2} \mu_i$. With probability at least $1 - 2\epsilon$, lil-KLUCB returns the arm with the largest mean and the total number of samples it collects is upper bounded by

$$X \leq c_0 \log \frac{1}{\epsilon} \log D(\mu, \mu_i) + c_0 \log(n+1) + \log D(\mu, \mu) \inf_{i \geq 2} \frac{1}{\mu_i - \mu} + \frac{1}{\epsilon} \log \frac{1}{\epsilon} \log D(\mu, \mu_i) + \frac{1}{\epsilon} \log D(\mu, \mu) \inf_{i \geq 2} \frac{1}{\mu_i - \mu}$$

where c_0 is some universal constant, $D(x, y)$ is the Chernoff-information.

Remark 1. Note that the LUCB++ algorithm of [11] is general enough to handle identification of the top k arms (not just the best-arm). All arguments presented in this paper also go through when considering the top- k problem for $k \geq 1$. However, to keep the arguments clear and concise, we chose to focus on the best-arm problem only. 2.1

Comparison with previous work

We now compare the sample complexity of lil-KLUCB to that of the two most closely related algorithms, KL-LUCB [8] and lil-UCB [5]. For a detailed review of the history of MAB problems and the use of KL-confidence intervals for bounded rewards, we refer the reader to [3, 9, 4]. For the KL-LUCB algorithm, Theorem 3 of [8] guarantees a high-probability sample complexity upper bound scaling as X

$$\inf_{i \geq 1} (D(\mu_i, c)) \log n + \log(D(\mu_i, c)) \cdot \frac{1}{c} \log \frac{1}{\epsilon}$$

Our result improves this in two ways. On one hand, we eliminate the unnecessary logarithmic dependence on the number of arms n in every term. Note that the $\log n$ factor still appears in Theorem 1 in the term corresponding to the number of samples on the best arm. It is shown in [11] that this factor is indeed unavoidable. The other improvement lil-KLUCB offers over KL-LUCB is improved logarithmic dependence on the Chernoff-information terms. This is due to the tighter confidence intervals derived in Section 3. Comparing Theorem 1 to the sample complexity of lil-UCB, we see that the two are of the same form, the exception being that the Chernoff-information terms take the place of the squared mean-gaps 3

(which arise due to the use of sub-Gaussian (SG) bounds). To give a sense of the improvement this can provide, we compare the sums $\sum_{i=1}^n \frac{1}{D(\mu_i, \mu)}$ and $\sum_{i=1}^n \frac{1}{D(\mu_i, \mu_i)}$

$$\sum_{i=1}^n \frac{1}{D(\mu_i, \mu)} \leq \sum_{i=1}^n \frac{1}{D(\mu_i, \mu_i)}$$

Let $\mu_i = (\mu_i, \sigma_i)$ and $\mu = (\mu, \sigma)$. Note that the Chernoff-information between $\text{Ber}(\mu_i)$ and $\text{Ber}(\mu)$ can be expressed as $D(\mu_i, \mu) = \max\{0, \min\{D(x, \mu_i), D(x, \mu)\}\} = D(x^*, \mu) = D(x^*, \mu_i) = x^*[\mu_i, \mu]$

$$D(x^*, \mu_i) + D(x^*, \mu) \geq 2$$

for some unique $x^* \in [\mu_i, \mu]$. It follows that $D(\mu_i, \mu) \geq \frac{1}{2}$

Consulting the proof of Theorem 1 it is clear that the number of samples on the sub-optimal arms of lil-KLUCB scales essentially as SKL w.h.p. (ignoring doubly logarithmic terms), and a similar argument can be made about lil-UCB. This justifies considering these sums in order to compare lil-KLUCB and lil-UCB.

4

(i) Define the sequence $z_t \in (0, 1]$, $t \leq N$ such that $\log(\frac{1}{N}) \log(2t)/\epsilon$
 $2D(\frac{1}{N} + \frac{1}{N^{b+1}} z_t, \epsilon) = t$ if a solution exists, and $z_t = 1$ otherwise.
Then $P(\exists t \leq N : \epsilon b t \leq \sum_{s=1}^t z_s) \leq \epsilon$.

(1)

(ii) Define the sequence $z_t \in [0, 1]$, $t \leq N$ such that
 $\log(\frac{1}{N}) \log(2t)/\epsilon \leq 2D(\frac{1}{N} + \frac{1}{N^{b+1}} z_t, \epsilon) = t$ if a solution exists, and $z_t = 1$ otherwise. Then $P(\exists t \leq N : \epsilon b t \leq \sum_{s=1}^t z_s) \leq \epsilon$. The result above can be used to construct anytime confidence bounds for the mean as follows. Consider part (i) of Theorem 2 and fix ϵ . The result gives a sequence z_t that upper bounds the deviations of the empirical mean. It is defined through an equation of the form $D(\frac{1}{N} + \frac{1}{N^{b+1}} z_t, \epsilon) = ft$. Note that the arguments of the function on the left must be in the interval $[0, 1]$, in particular $\frac{1}{N^{b+1}} z_t \leq 1$, and the maximum of $D(\frac{1}{N} + x, \epsilon)$ for $x \in [0, 1]$ is $D(1, \epsilon) = \log \frac{1}{\epsilon}$. Hence, equation 1 does not have a solution if t is too large (that is, if t is small). In these cases we set $z_t = 1$. However, since ft is decreasing, equation 1 does have a solution when $t \geq T$ (for some T depending on ϵ), and this solution is unique (since $D(\frac{1}{N} + x, \epsilon)$ is strictly increasing). With high probability $\epsilon b t \leq \sum_{s=1}^t z_s$ for all $t \leq N$ by Theorem 2. Furthermore, the function $D(\frac{1}{N} + x, \epsilon)$ is increasing in $x \geq 0$. By combining these facts we get that with probability at least $1 - \epsilon$

$t \leq \frac{1}{\epsilon} D(\frac{1}{N} + \frac{1}{N^{b+1}} z_t, \epsilon) \leq \frac{1}{\epsilon} D(\frac{1}{N^{b+1}}, \epsilon)$. On the other hand

$\log(\frac{1}{N}) \log(2t)/\epsilon \leq t$ by definition. Chaining these two inequalities leads to the lower confidence bound

$$\log(\frac{1}{N}) \log(2t)/\epsilon \leq N \epsilon b t + m \leq 2L(t, \epsilon) = \inf_{m \leq \epsilon b t} D(\frac{1}{N^{b+1}}, m) +$$

$$D(\frac{1}{N^{b+1}} z_t, \epsilon)$$

?

(2)

which holds for all times t with probability at least $1 - \epsilon$. Considering the left deviations of $\epsilon b t \leq \sum_{s=1}^t z_s$ we can get an upper confidence bound in a similar manner:

$$\log(\frac{1}{N}) \log(2t)/\epsilon + m \leq 2U(t, \epsilon) = \sup_{m \leq \epsilon b t} D(\frac{1}{N^{b+1}}, m) +$$

(3) t That is, for all times t , with probability at least $1 - \epsilon$ we have $L(t, \epsilon) \leq \epsilon b t \leq U(t, \epsilon)$. Note that the constant $\log(\frac{1}{N}) \leq 2 \log(2N)$, so the choice of N plays a relatively mild role in the bounds. However, we note here that if N is sufficiently large, then $\frac{1}{N^{b+1}} \leq \epsilon b t$, and thus

$N \epsilon b t + m \leq D(\frac{1}{N^{b+1}}, m) \leq D(\epsilon b t, m)$, in which case the bounds above are easily compared to those in prior works [4, 9, 3]. We make this connection more precise and show that the confidence intervals defined as

$$c(N) \log(\frac{1}{N}) \log(2t)/\epsilon \leq L_0(t, \epsilon) = \inf_{m \leq \epsilon b t} D(\epsilon b t, m) \leq \epsilon b t$$

j?[N]

$\sum_{j=1}^N \sum_{k=1}^N (N+1) D_j + \sum_{j=1}^N z_{tj}^j, \sum_{j=1}^N D_j + N+1 \sum_{j=1}^N z_{tj}^j, \dots$. This implies $P(A_k) = 0$.

•

Plugg

Regarding the first term in (4), again using the Bernoulli rate function bound we have

.

i?[N]

8

Analysis of lil-KLUCB

Recall that the lil-KLUCB algorithm uses confidence bounds of the form $U_i(t, \beta) = \sup\{m_i : b_i(t, m_i) \leq f_i(t, \beta)\}$ with some decreasing sequence $f_i(t, \beta)$. In this section we make this dependence explicit, and use the notations $U_i(f_i(t, \beta))$ and $L_i(f_i(t, \beta))$ for upper and lower confidence bounds. For any $i \in [0]$ and $t \in [n]$, define the events $\mathcal{E}_i(t) = \{t \in N : L_i(f_i(t, \beta)), U_i(f_i(t, \beta))\}$. The correctness of the algorithm follows from the correctness of the individual confidence intervals, as is usually the case with LUCB algorithms. This is shown formally in Proposition 1 provided in the Supplementary Materials. The main focus in this section is to show a high probability upper bound on the sample complexity. This can be done by combining arguments frequently used for analyzing LUCB algorithms and those used in the analysis of the lil-UCB [5]. The proof is very similar in spirit to that of the LUCB++ algorithm [11]. Due to spatial restrictions, we only provide a proof sketch here, while the detailed proof is provided in the Supplementary Materials. Proof sketch of Theorem 1. Observe that at each time step two things can happen (apart from stopping): (1) Arm 1 is not sampled (two sub-optimal arms are sampled); (2) Arm 1 is sampled together with some other (suboptimal) arm. Our aim is to upper bound the number of times any given arm is sampled for either of the reasons above. We do so by conditioning on the event \mathcal{E}_1 .

$\mathcal{E}_0 = \mathcal{E}_1$ for a certain choice of $\{\mathcal{E}_i\}$ defined below.

For instance, if arm 1 is not sampled at a given time t , we know that $\text{TOP}(t) = 1$, which means there must be an arm $i \in [2]$ such that $U_i(T_i(t), \beta) \leq U_1(T_1(t), \beta)$. However, on the event \mathcal{E}_1 , the UCB of arm 1 is accurate, implying that $U_1(T_1(t), \beta) \leq \beta_1$. This implies that $T_i(t)$ can not be too big, since on \mathcal{E}_i , β_i is "close" to β_i , and also $U_i(T_i(t), \beta)$ is not much larger than β_i . All this is made formal in Lemma 2, yielding the following upper bound on number of times arm i is sampled for reason (1): $\sum_{t \in N} \mathbb{1}_{\mathcal{E}_i} = \min\{t \in N : f_i(t, \beta) \leq D(\beta_i, \beta_1)\}$. Similar arguments can be made about the number of samples of any suboptimal arm i for reason (2), and also the number of samples on arm 1. This results in the sample complexity upper bound

$$\sum_{i=1}^K K_1 \log \frac{1}{\beta_1} \log D(\beta_i, \beta_1) + \log \frac{1}{\beta_1} K_1 \log(n \beta_1) + \log D(\beta_1, \beta_e) + D(\beta_1, \beta_e) D(\beta_i, \beta_e) \frac{1}{\beta_1^2}$$

0

on the event \mathcal{E}_1 , where K_1 is a universal constant. Finally, we define the quantities $\mathcal{E}_i = \sup\{t \in N : U_i(f_i(t, \beta)) \leq \beta_i\}$. Note that we have $P(\mathcal{E}_i) = P(t \in N : U_i(f_i(t, \beta)) \leq \beta_i)$ according to Theorem 1 in the Supplementary Material. Substituting $\beta_i = \exp(D(\beta_i, \beta_e)z)$ we get

$\log \frac{1}{\beta_1} \log D(\beta_i, \beta_e) = z \exp(D(\beta_i, \beta_e)z)$. Hence $\{\mathcal{E}_i\}$ are independent sub-exponential variables, which allows us to control their contribution to the sum above using standard techniques.

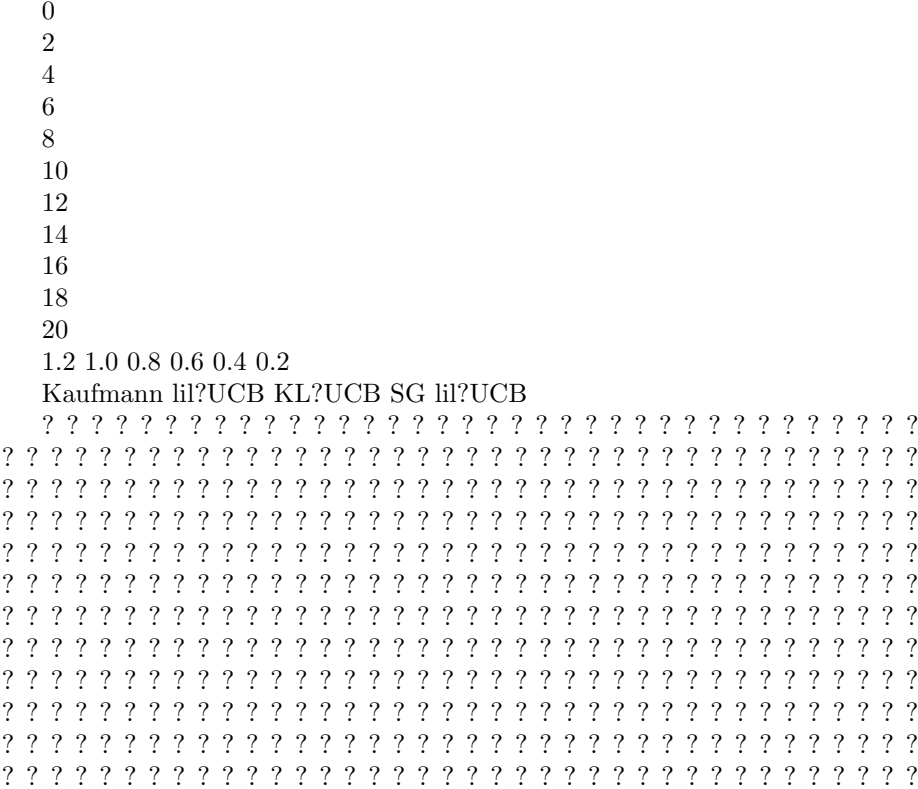
Real-World Crowdsourcing

We now compare the performance of lil-KLUCB to that of other algorithms in the literature. We do this using both synthetic data and real data from the

Number of samples (10 thousands)
0.15
0.30
0.45
0.60
0.75
0.90
Number of samples (10 thousands)

Figure 1: Probability of the best-arm in the top 5 empirically best arms, as a function of the number of samples, based on 250 repetitions. $\mu_i = 1 - ((i - 1)/n)^\alpha$, with $\alpha = 1$ in the left panel, and $\alpha = 1/2$ in the right panel. The mean-profile is shown above each plot. [KL] Blue; [SG1] Red; [SG2] Black.

As seen in Table 1, the KL confidence bounds have the potential to greatly outperform the subGaussian ones. To illustrate this indeed translates into superior performance, we simulate two cases, with means $\mu_i = 1 - ((i - 1)/n)^\alpha$, with $\alpha = 1/2$ and $\alpha = 1$, and $n = 1000$. As expected, the KL-based method requires significantly fewer samples (about 20 % for $\alpha = 1$ and 30 % for $\alpha = 1/2$) to find the best arm. Furthermore, the arms with means below the median are sampled about 15 and 25 % of the time respectively – key in crowdsourcing applications, since having participants answer fewer irrelevant (and potentially annoying) questions improves both efficiency and user experience.



?
 ?
 ?
 ?
 ?

0.0
 (Empirical) probability ? 250 trials
 P(best arm in top 5), Contest 558
 0.00
 Number of samples (10 thousands)
 0.75
 1.50
 2.25
 3.00
 3.75
 4.50
 Number of samples (10 thousands)

Figure 2: Probability of the best-arm in the top 5 empirically best arms vs. number of samples, based on 250 bootstrapped repetitions. Data from New Yorker contest 558 ($\mu_1 = 0.536$) on left, and contest 512 ($\mu_1 = 0.8$) on right. Mean-profile above each plot. [KL] Blue; [SG1] Red; [SG2] Black.

8

To see how these methods fair on real data, we also run these algorithms on bootstrapped human response data from the real New Yorker Caption Contest. The mean reward of the best arm in these contests is usually between 0.5 and 0.85, hence we choose one contest from each end of this spectrum. At the lower end of the spectrum, the three methods fair comparably. This is expected because the sub-Gaussian bounds are relatively good for means about 0.5. However, in cases where the top mean is significantly larger than 0.5 we see a marked improvement in the KL-based algorithm.

Extension to numerical experiments Since a large number of algorithms have been proposed in the literature for best arm identification, we include another algorithm in the numerical experiments for comparison. Previously we compared lil-KLUCB to lil-UCB as a comparison for two reasons. First, this comparison illustrates best the gains of using the novel anytime confidence bounds as opposed to those using sub-Gaussian tails. Second, since lil-UCB is the state of the art algorithm, any other algorithm will likely perform worse. The authors of [6] compare a number of different best arm identification methods, and conclude that two of them seem to stand out: lil-UCB and Thompson sampling. Therefore, we now include Thompson sampling [Th] in our numerical experiments for the New Yorker data. We implemented the method as prescribed in [6]. As can be seen in Figure 5, Thompson sampling seems to perform somewhat worse than the previous methods in these two instances.

0
 2
 4

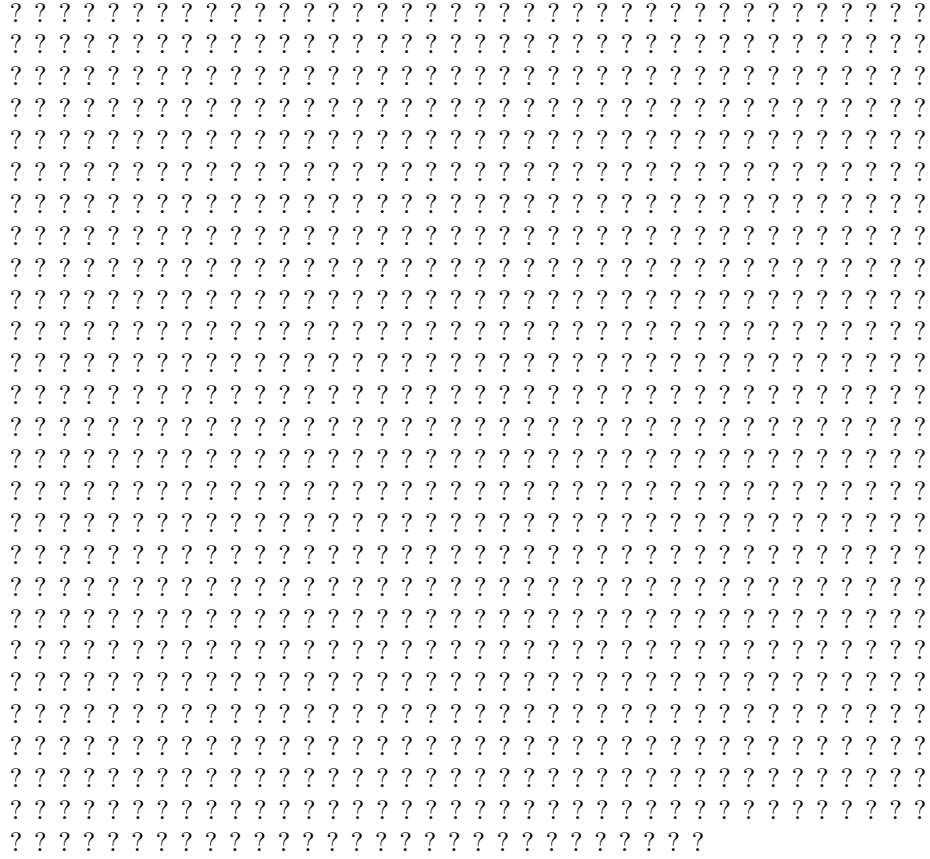


Figure 3: Probability of the best-arm in the top 5 empirically best arms vs. number of samples, based on 250 bootstrapped repetitions. Data from New Yorker contest 558 ($p_1 = 0.536$) on left, and contest 512 ($p_1 = 0.8$) on right. Mean-profile above each plot. [KL] Blue; [SG1] Red; [SG2] Black; [Th] Purple.

2 References

- [1] Jean-Yves Audibert and Sébastien Bubeck. Best arm identification in multi-armed bandits. In COLT-23th Conference on Learning Theory-2010, pages 13?p, 2010. [2] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. Concentration inequalities: A nonasymptotic theory of independence. Oxford university press, 2013. [3] Olivier Cappé, Aurélien Garivier, Odalric-Ambrym Maillard, Rémi Munos, Gilles Stoltz, et al. Kullback–leibler upper confidence bounds for optimal sequential allocation. *The Annals of Statistics*, 41(3):1516?1541, 2013. [4] Aurélien Garivier and Olivier Cappé. The kl-ucb algorithm for bounded stochastic bandits and beyond. In COLT, pages 359?376, 2011. [5] Kevin Jamieson, Matthew Malloy, Robert Nowak, and Sébastien Bubeck. lil-ucb: An optimal exploration algorithm for multi-armed bandits. In Conference on Learning Theory, pages 423?439, 2014. [6] Kevin G Jamieson, Lalit Jain, Chris Fernandez, Nicholas J Glattard, and Rob Nowak. Next: A system for real-world development, evaluation, and application of active learning. In *Advances in Neural Information Processing Systems*, pages 2656?2664, 2015. [7] Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best arm identification in multi-armed bandit models. *The Journal of Machine Learning Research*, 2016. [8] Emilie Kaufmann and Shivaram Kalyanakrishnan. Information complexity in bandit subset selection. In COLT, pages 228?251, 2013. [9] Odalric-Ambrym Maillard, Rémi Munos, Gilles Stoltz, et al. A finite-time analysis of multi-armed bandits problems with kullback-leibler divergences. In COLT, pages 497?514, 2011. [10] B. Fox Rubin. How new yorker cartoons could teach computers to be funny. CNET News, 2016. <https://www.cnet.com/news/how-new-yorker-cartoons-could-teach-computers-to-be-funny/>. [11] Max Simchowitz, Kevin Jamieson, and Benjamin Recht. The simulator: Understanding adaptive sampling in the moderate-confidence regime. arXiv preprint arXiv:1702.05186, 2017.