

# Corrupt Bandits for Preserving Local Privacy

---

Pratik Gajane<sup>1</sup> Tanguy Urvoy<sup>2</sup> Emilie Kaufmann<sup>3</sup>

7<sup>th</sup> April 2018

<sup>1</sup>Montanuniversität Leoben

<sup>2</sup>Orange labs

<sup>3</sup>CNRS & Univ. Lille & Inria-Sequel

Motivation and Formalization

Lower Bound on Regret

Proposed Algorithms

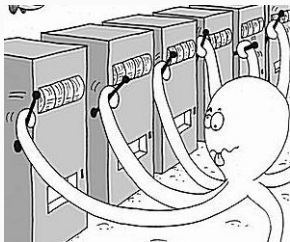
Experiments

Final Remarks

# Motivation and Formalization

---

# Classical Stochastic Bandits



- $K$  arms/actions
- Unknown reward distributions with mean  $\mu_a$  for arm  $a$
- Learner pulls arm  $a$ 
  - receives reward  $\sim$  distribution for  $a$
  - feedback = received reward  
(**Absolute feedback**)
- **Regret** = best possible reward - reward of pulled arm
- Learner's goal = minimize **cumulative regret**

## Motivation for Corrupt Bandits: Privacy

---

# Motivation for Corrupt Bandits: Privacy



"If you're doing something that you don't want other people to know, maybe you shouldn't be doing it in first place"



"Privacy is no longer a social norm!"

# Local Differential Privacy

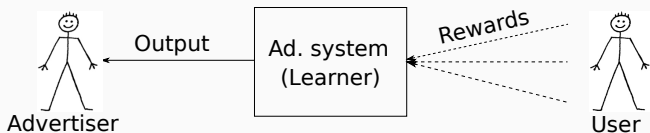


Figure 1: Ad system using bandits

- Ad application as bandit problem.
- Feedback from users on ads (arms).
- Information about user tastes as output to advertisers.

# Local Differential Privacy

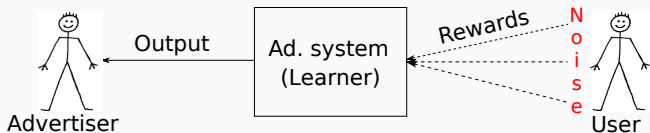


Figure 1: Ad system using bandits

- Ad application as bandit problem.
- Feedback from users on ads (arms).
- Information about user tastes as output to advertisers.
- Local differential privacy (DP), by Duchi et al.(2014) [2].
- Classical bandits unable to deal with noisy feedback.



# Questions???

- Bandit setting to deal with Corrupted/Noisy Feedback?
- Regret Lower Bound for such Bandit setting?
- Algorithms to solve this Bandit setting?

# Corrupt Bandits: Formalization

- Formally characterized by
  - $K$  arms
  - unknown reward distribution with mean  $\mu_a$  for each  $a$
  - unknown feedback distribution with mean  $\lambda_a$  for each  $a$
  - known mean corruption function  $g_a$  for each  $a$
- $g_a(\mu_a) = \lambda_a$
- Learner's goal: minimize cumulative **regret**

## Lower Bound on Regret

---

Theorem (Thm. 1, PG, Urvoy & Kaufmann(2018) [4])

Any algorithm for a Bernoulli corrupt bandit problem satisfies,

$$\liminf_{T \rightarrow \infty} \frac{\text{Regret}_T}{\log(T)} \geq \sum_{a=2}^K \frac{\Delta_a}{d(\lambda_a, g_a(\mu_1))}.$$

where  $d(x, y) := \text{KL}(\mathcal{B}(x), \mathcal{B}(y))$

- $\Delta_a$  = optimal mean reward - mean reward of a ( $\mu_a$ )
- 1 is assumed to be the optimal arm w.l.o.g.
- $\lambda_a = g_a(\mu_a)$ . Behaviour of  $g_a$  on  $\mu_a$  and  $\mu_1$  affects lower bound.

# Proposed Algorithms

---

## Proposed algorithm: kl-UCB-CF

### Algorithm: kl-UCB-CF

*Pull at time  $t$  an arm maximizing*

$$\text{Index}_a(t) := \max\{q : N_a(t) \cdot d(\hat{\lambda}_a(t), g_a(q)) \leq f(t)\}$$

- Similar to kl-UCB by Cappé et al. (2013) [1] for classical bandits.
- $\text{Index}_a(t)$  = UCB on  $\mu_a$  from confidence interval on  $\lambda_a$  and using exploration function  $f$
- $\hat{\lambda}_a(t)$  = emp. mean of feedback of  $a$  until time  $t$
- UCB1 (Auer et al. (2002)) can be updated to UCB-CF.

# Upper Bound for $\text{kl-UCB-CF}$

Theorem (Thm. 2, PG, Urvoy & Kaufmann(2018) [4])

$$\text{Regret of kl-UCB-CF} \leq \sum_{a=2}^K \frac{\Delta_a \log(T)}{d(\lambda_a, g_a(\mu_1))} + O(\sqrt{\log(T)}).$$

- Recall that 1 is assumed to be the optimal arm.
- More explicit bound can be provided.
- Optimal as upper bound matches lower bound.

# Proposed Algorithm: TS-CF

## Algorithm: TS-CF

1. Sample  $\theta_a(t)$  from Beta posterior distribution on mean feedback of arm  $a$ .
2. Pull arm  $\hat{a}_{t+1} = \arg \max_a g_a^{-1}(\theta_a(t))$ .

- Similar to Thompson sampling by Thompson (1933) [5] for classical bandits.
- Probability ( $a$  is played) = posterior probability ( $a$  is optimal).



Theorem (Thm. 3, PG, Urvoy & Kaufmann(2018) [4])

$$\text{Regret of TS-CF} \leq \sum_{a=2}^K \frac{2\Delta_a \log(T)}{d(\lambda_a, g_a(\mu_1))} + O(\sqrt{\log(T)})$$

- Recall that 1 is assumed to be the optimal arm.
- A tighter bound can be provided.
- Optimal as upper bound matches lower bound.

# Experiments

---

# Experiments with varying time

- Bernoulli corrupt bandit:  $\mu_1 = 0.9$        $\mu_2 = \dots = \mu_{10} = 0.6$
- Comparison over a period of time for fixed corruption

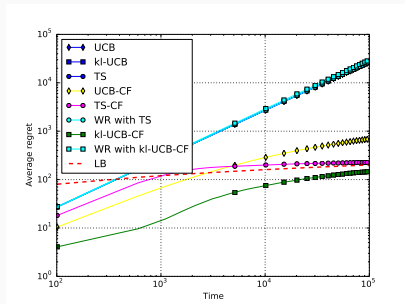


Figure 2: Regret plots with varying  $T$  up to  $10^5$

# Experiments with varying Local DP

- Bernoulli corrupt bandit:  $\mu_1 = 0.9$        $\mu_2 = \dots = \mu_{10} = 0.6$
- Comparison with varying level of Local DP;  $\epsilon$  from  $\{1/8, 1/4, 1/2, 1, 2, 4, 8\}$

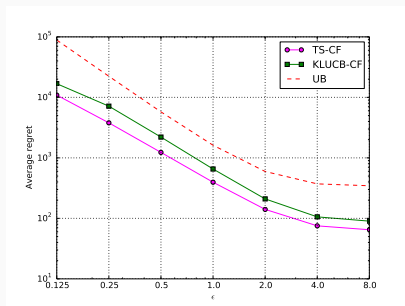


Figure 3: Regret with varying level of Local DP

## Final Remarks

---

Covered in this talk:

- Introduced Corrupt Bandits to provide privacy.
- Proved the lower bound. Provided optimal algorithms matching the lower bound.

Covered in this talk:

- Introduced Corrupt Bandits to provide privacy.
- Proved the lower bound. Provided optimal algorithms matching the lower bound.

Not covered in this talk:

- Provided optimal mechanism for achieving local DP.

Covered in this talk:

- Introduced Corrupt Bandits to provide privacy.
- Proved the lower bound. Provided optimal algorithms matching the lower bound.

Not covered in this talk:

- Provided optimal mechanism for achieving local DP.
- Proved regret guarantees for achieving required level of local DP (Trade-off between utility and privacy).



# Final Remarks

Covered in this talk:

- Introduced Corrupt Bandits to provide privacy.
- Proved the lower bound. Provided optimal algorithms matching the lower bound.

Not covered in this talk:

- Provided optimal mechanism for achieving local DP.
- Proved regret guarantees for achieving required level of local DP (Trade-off between utility and privacy).

Future work:

- Contextual corruption?

# Final Remarks

Covered in this talk:

- Introduced Corrupt Bandits to provide privacy.
- Proved the lower bound. Provided optimal algorithms matching the lower bound.

Not covered in this talk:

- Provided optimal mechanism for achieving local DP.
- Proved regret guarantees for achieving required level of local DP (Trade-off between utility and privacy).

Future work:

- Contextual corruption?
- Corrupted feedback in RL? (a recent publication by Everitt et al. (2017) [3]).

Thank you all.

## References

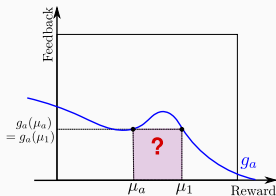
---

- [1] Olivier Cappé, Aurélien Garivier, Odalric-Ambrym Maillard, Rémi Munos, and Gilles Stoltz. Kullback-Leibler upper confidence bounds for optimal sequential allocation. *Annals of Statistics*, 41(3):1516–1541, 2013.
- [2] John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. Privacy aware learning. *J. ACM*, 61(6):38:1–38:57, December 2014.
- [3] Tom Everitt, Victoria Krakovna, Laurent Orseau, and Shane Legg. Reinforcement learning with a corrupted reward channel. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, pages 4705–4713, 2017.

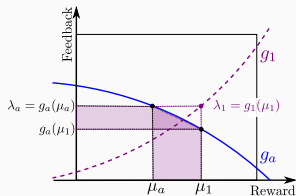
- [4] Pratik Gajane, Tanguy Urvoy, and Emilie Kaufmann. Corrupt bandits for preserving local privacy. In *Proceedings of the 29th International Conference on Algorithmic Learning Theory (ALT)*, 2018.
- [5] W.R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Bulletin of the AMS*, 25:285–294, 1933.

# Interpretation of Lower Bound for Corrupt Bandits

- Divergence between  $\lambda_a$  and  $g_a(\mu_1)$  plays a crucial role in distinguishing arm  $a$  from the optimal arm.



(a) Uninformative  $g_a$  function.



(b) Informative  $g_a$  function.

**Figure 4:** On the left,  $g_a$  is such that  $\lambda_a = g_a(\mu_1)$ . On the right, a steep monotonic  $g_a$  leads  $\Delta_a = \mu_1 - \mu_a$  into a clear gap between  $\lambda_a$  and  $g_a(\mu_1)$ .

- If the  $g_a$  function is non-monotonic, it might be impossible to distinguish between arm  $a$  and the optimal arm.
- Assumption: Corruption functions strictly monotonic.

# Optimal mechanism for local DP and regret

- Corruption matrix

$$\mathbb{M}_a = \begin{matrix} & 0 & 1 \\ \begin{matrix} 0 \\ 1 \end{matrix} & \begin{bmatrix} \frac{e^\epsilon}{1+e^\epsilon} & \frac{1}{1+e^\epsilon} \\ \frac{1}{1+e^\epsilon} & \frac{e^\epsilon}{1+e^\epsilon} \end{bmatrix} \end{matrix}.$$

## Corollary

*The regret of kl-UCB-CF or TS-CF at time  $T$  with  $\epsilon$ -locally differentially private bandit feedback corruption scheme is*

$$\text{Regret}_T \leq \sum_{a=2}^K \frac{2 \log(T)}{\Delta_a \left(\frac{e^\epsilon - 1}{e^\epsilon + 1}\right)^2} + O(\sqrt{\log(T)}).$$

# Local DP vs global DP

- For low values of  $\epsilon$ ,  $\left(\frac{e^\epsilon - 1}{e^\epsilon + 1}\right) \approx \epsilon/2$ .
- In-line with global DP algorithms with a multiplicative factor of  $O(\epsilon^{-1})$  or  $O(\epsilon^{-2})$ .
- One global DP algorithm with additive factor of  $O(\epsilon^{-1})$ . Our lower bound shows that's not possible for local DP.

