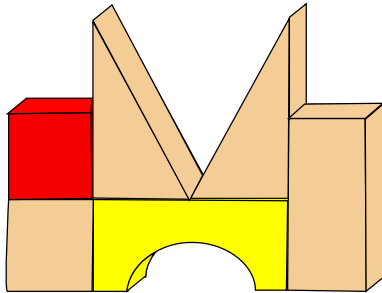


Matching Regret Lower Bounds in Structured Stochastic Bandits



Wouter M. Koolen

CWI

Centrum Wiskunde & Informatica

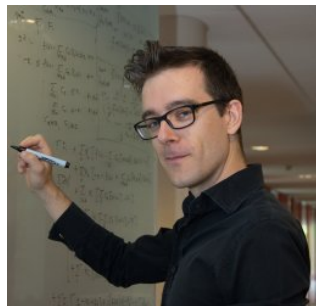
Team



Rémy Degenne

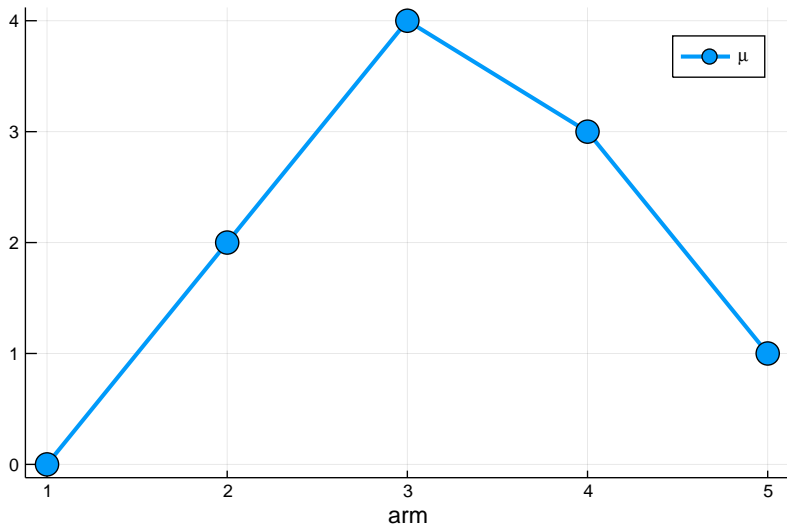


Han Shao (邵涵)

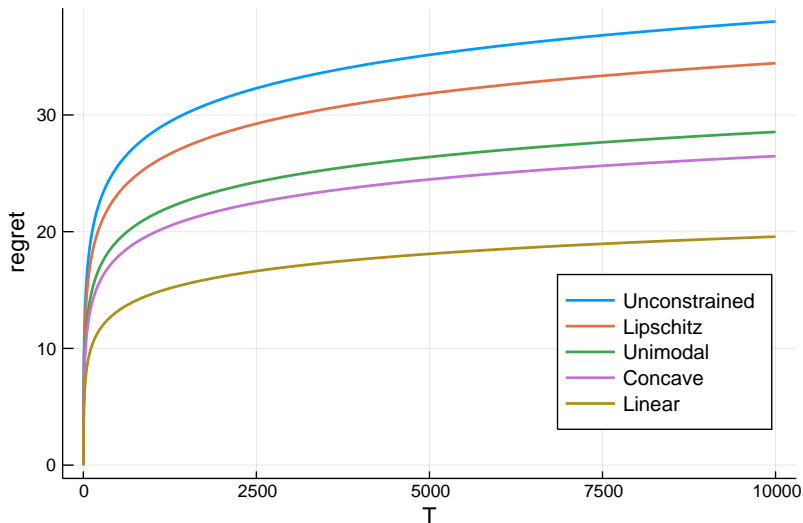
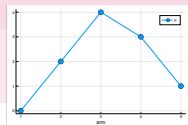


Wouter Koolen

Stochastic Bandit Instance (Running Example)



Desired behaviour

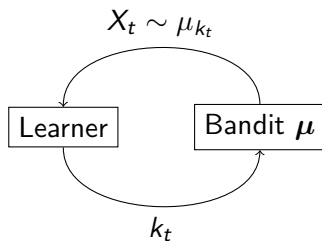




Outline

- 1 Introduction
- 2 Lower bound
- 3 Noise Free Case
- 4 The Real Deal
- 5 Experiments

Setting



Structure $\mathcal{M} \subseteq R^K$.

MAB instance $\mu \in \mathcal{M}$

Expfam $d(\mu, \lambda)$

Gaps $\Delta^k = \mu^* - \mu^k$

Regret

$$\sum_{t=1}^T \mathbb{E}[\Delta^{k_t}]$$

Goals



- Asymptotic Optimality
- Finite-time Regret Guarantees
- General Structure-Aware Methodology
- Computational Efficiency

Banditual Context

Regret

- Unimodal [Combes and Proutiere, 2014]
- Lipschitz [Magureanu, Combes, and Proutière, 2014]
- Rank-1 [Katariya, Kveton, Szepesvári, Vernade, and Wen, 2017]
- Linear [Lattimore and Szepesvári, 2017]
- OSSB [Combes, Magureanu, and Proutiere, 2017]

Pure Exploration

- Track-and-Stop (MAB) [Garivier and Kaufmann, 2016]
- Structure, Gaussian [Chen, Gupta, Li, Qiao, and Wang, 2017]
- Structure, ExpFam [Kaufmann and Koolen, 2018]
- Game core [Degenne, Koolen, and Ménard, 2019] **yesterday**



Outline

- 1 Introduction
- 2 Lower bound**
- 3 Noise Free Case
- 4 The Real Deal
- 5 Experiments

Argument [Graves and Lai, 1997]

Fix an **asymptotically consistent** algorithm for structure \mathcal{M} . Consider its behaviour on $\mu \in \mathcal{M}$, and on any alternative bandit model $\lambda \in \mathcal{M}$ with $i^*(\mu) \neq i^*(\lambda)$:

$$\mathbb{E}_{\mu}[N_T^{i^*(\mu)}]/T \rightarrow 1 \quad \text{but} \quad \mathbb{E}_{\lambda}[N_T^{i^*(\mu)}]/T \rightarrow 0.$$

This stark **difference in behaviour** requires **discriminating information!**
Specifically,

$$\text{KL}(\mathbb{P}_{\mu}^T \parallel \mathbb{P}_{\lambda}^T) = \sum_k \mathbb{E}_{\mu}[N_T^k] d(\mu^k, \lambda^k) \geq \ln T.$$

Instance-Dependent Regret Lower Bound



Any asymptotically consistent algorithm for structure \mathcal{M} must incur on each $\mu \in \mathcal{M}$ regret at least

$$V_T = \min_{N \geq 0} \sum_k N^k \Delta^k \quad \text{subject to} \quad \inf_{\lambda \in \Lambda} \sum_k N^k d(\mu^k, \lambda^k) \geq \ln T$$

where

$$\Lambda = \{\lambda \in \mathcal{M} \mid i^*(\lambda) \neq i^*(\mu)\}$$

This is a (semi-infinite) **covering linear program**.

Operationalising the Lower Bound

Earlier work

At each time step

- compute **oracle sample counts** $N^*(\hat{\mu}_t)$ and advance $N_t \rightarrow N^*$, or
- **force exploration** to ensure $\hat{\mu}_t \rightarrow \mu$.

Operationalising the Lower Bound

Earlier work

At each time step

- compute **oracle sample counts** $N^*(\hat{\mu}_t)$ and advance $N_t \rightarrow N^*$, or
- **force exploration** to ensure $\hat{\mu}_t \rightarrow \mu$.

This talk

- Reformat lower bound as zero-sum “minigame”.
- **Iteratively** solve minigame by full information online learning.
- Use iterates to advance N_t .
- Add optimism to induce exploration.
- **Compose** regret bound from minigame regret + estimation regret

Minigame

We have $V_T = \frac{\ln T}{D^*}$ where

$$D^* = \max_{w \in \Delta} \inf_{\lambda \in \Lambda} \frac{\sum_k w^k d(\mu^k, \lambda^k)}{\sum_k w^k \Delta^k}$$

$w^k \propto N^k$
pulls

Minigame

We have $V_T = \frac{\ln T}{D^*}$ where

$$D^* = \max_{w \in \Delta} \inf_{\lambda \in \Lambda} \frac{\sum_k w^k d(\mu^k, \lambda^k)}{\sum_k w^k \Delta^k}$$

$w^k \propto N^k$
pulls

$$= \max_{\tilde{w} \in \Delta} \inf_{\lambda \in \Lambda} \sum_k \tilde{w}^k \frac{d(\mu^k, \lambda^k)}{\Delta^k}$$

$\tilde{w}^k \propto N^k \Delta^k$
regret

Minigame

We have $V_T = \frac{\ln T}{D^*}$ where

$$D^* = \max_{w \in \Delta} \inf_{\lambda \in \Lambda} \frac{\sum_k w^k d(\mu^k, \lambda^k)}{\sum_k w^k \Delta^k}$$

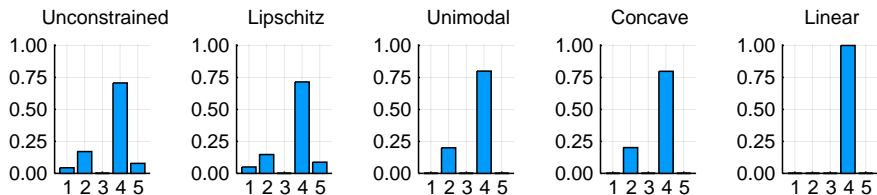
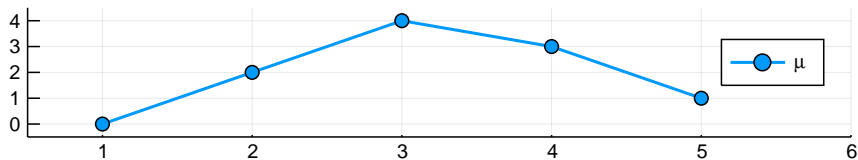
$w^k \propto N^k$
pulls

$$= \max_{\tilde{w} \in \Delta} \inf_{\lambda \in \Lambda} \sum_k \tilde{w}^k \frac{d(\mu^k, \lambda^k)}{\Delta^k}$$

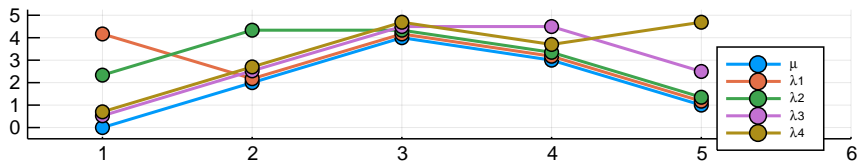
$\tilde{w}^k \propto N^k \Delta^k$
regret

$$= \inf_{q \in \Delta(\Lambda)} \max_k \frac{\mathbb{E}_{\lambda \sim q} [d(\mu^k, \lambda^k)]}{\Delta^k}$$

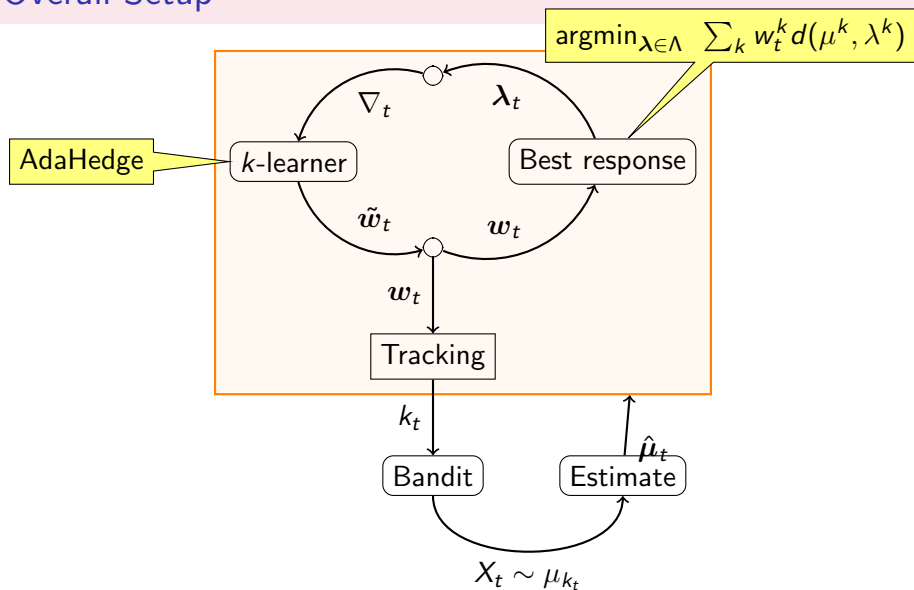
Illustration



Support for Lipschitz



Overall Setup





Outline

- 1 Introduction
- 2 Lower bound
- 3 Noise Free Case**
- 4 The Real Deal
- 5 Experiments



Noise-free result

Let \mathcal{B}_n^k be regret of full information online learning (AdaHedge) w. linear losses on the simplex.

Theorem

Consider running our algorithm until $\inf_{\lambda \in \Lambda} \sum_{t=1}^n \sum_k w_t^k d(\mu^k, \lambda^k) \geq \ln T$.
The iterates w_1, \dots, w_n satisfy

$$R_n = \sum_{t=1}^n \langle w_t, \Delta \rangle \leq V_T + \frac{\mathcal{B}_n^k}{D^*}$$

Note

- Can get k_1, \dots, k_n using tracking (at cost $\Delta^{\max} \ln K$)
- Standard choice gives $n = O(\ln T)$ and $\mathcal{B}_n^k = O(\sqrt{n}) = O(\sqrt{\ln T}) = o(\ln T)$.



Regret analysis

Given moves $\mathbf{w}_t \in \Delta_K$ and $\lambda_t \in \Lambda$, we instantiate a k -learner for the gain function

$$g_t(\tilde{\mathbf{w}}) = \langle \mathbf{w}_t, \Delta \rangle \sum_k \tilde{\mathbf{w}}^k \frac{d(\mu^k, \lambda_t^k)}{\Delta^k}$$

to provide regret bound

$$\sum_{t=1}^n g_t(\tilde{\mathbf{w}}_t) \geq \max_k \sum_{t=1}^n \langle \mathbf{w}_t, \Delta \rangle \frac{d(\mu^k, \lambda_t^k)}{\Delta^k} - \mathcal{B}_n^k. \quad (1)$$



Regret analysis (ctd)

Given \tilde{w}_t from the k -learner, we define player and opponent by

$$w_t^k \propto \tilde{w}_t^k / \Delta^k \quad (2)$$

$$\lambda_t \in \operatorname{argmin}_{\lambda \in \Lambda} \sum_k w_t^k d(\mu^k, \lambda^k) \quad (3)$$

to obtain

$$\begin{aligned} \sum_{t=1}^n g_t(\tilde{w}_t) &= \sum_{t=1}^n \langle w_t, \Delta \rangle \sum_k \tilde{w}_t^k \frac{d(\mu^k, \lambda_t^k)}{\Delta^k} \stackrel{(2)}{=} \sum_{t=1}^n \sum_k w_t^k d(\mu^k, \lambda_t^k) \\ &\stackrel{(3)}{=} \sum_{t=1}^n \inf_{\lambda \in \Lambda} \sum_k w_t^k d(\mu^k, \lambda^k) \leq \inf_{\lambda \in \Lambda} \sum_{t=1}^n \sum_k w_t^k d(\mu^k, \lambda^k) \end{aligned} \quad (4)$$



Regret analysis (ctd)

The stopping condition plus regret bounds (1) and (4) result in

$$\begin{aligned} \ln T + \mathcal{B}_n^k &\geq \max_k \sum_{t=1}^n \langle \mathbf{w}_t, \Delta \rangle \frac{d(\mu^k, \lambda_t^k)}{\Delta^k} = R_n \max_k \sum_{t=1}^n \frac{\langle \mathbf{w}_t, \Delta \rangle}{R_n} \frac{d(\mu^k, \lambda_t^k)}{\Delta^k} \\ &\geq R_n \inf_{q \in \Delta(\Lambda)} \max_k \frac{\mathbb{E}_{\lambda \sim q} [d(\mu^k, \lambda^k)]}{\Delta^k} = R_n D^* \end{aligned}$$

where we abbreviated $R_n = \sum_{t=1}^n \langle \mathbf{w}_t, \Delta \rangle$. All in all we showed

$$R_n \leq V_T + \frac{\mathcal{B}_n^k}{D^*}$$

On Symmetry

Game-theoretic equilibrium is **symmetric** concept.

Can also focus on λ -learner instead of k -learner. Interesting trade-offs

- More complex domain $\lambda \in \Lambda$.
- No need for tracking, best response in k is “pure” arm.

Will show both in experiments.



Outline

- 1 Introduction
- 2 Lower bound
- 3 Noise Free Case
- 4 The Real Deal**
- 5 Experiments

Scaling up

Can use what we developed so far to compute oracle weights every round (OSSB). Efficient for **every** bandit structure for which best response is tractable.

Scaling up

Can use what we developed so far to compute oracle weights every round (OSSB). Efficient for **every** bandit structure for which best response is tractable.

But we can do much better!

Scaling up

Can use what we developed so far to compute oracle weights every round (OSSB). Efficient for **every** bandit structure for which best response is tractable.

But we can do much better!

Idea:

- Run only one iteration every round.
- Deal with unknown μ .
- Exploitation.

some issues . . .

First Issue

Actually, $\Delta^* = 0$. And we were dividing by it all over the place.

First Issue

Actually, $\Delta^* = 0$. And we were dividing by it all over the place.

Idea: run on $\Delta_\epsilon^k = \max\{\Delta^k, \epsilon\}$.

Theorem

$$\lim_{\epsilon \rightarrow 0} V_T^\epsilon = V_T$$

In several cases we can show perturbed value is $V_T^\epsilon \leq V_T + \sqrt{2\epsilon V_T}$.

One iteration every round

- Replace μ by **estimate** $\hat{\mu}_t$.
- Add **optimism** to force exploration.

We introduce upper confidence bounds on the ratio KL/gap.

$$\text{UCB}_s^k = \sup_{\xi \in \mathcal{C}_{s-1}^k} \frac{d(\xi, \lambda_t^k)}{\max \{ \epsilon_s, \mathbf{1}\{k \neq j_s\} [\mu_{s-1}^+ - \xi] \}}$$

$$\text{where } \mathcal{C}_{s-1}^k = \left[\hat{\mu}_{s-1}^k \pm \sqrt{\frac{\overline{\ln}(n_{s-1}^j, N_{s-1}^k)}{N_{s-1}^k}} \right].$$

- We do not know **identity of the best arm**, and hence Λ (domain of λ) Estimate best arm, and run K independent interactions.

Algorithm

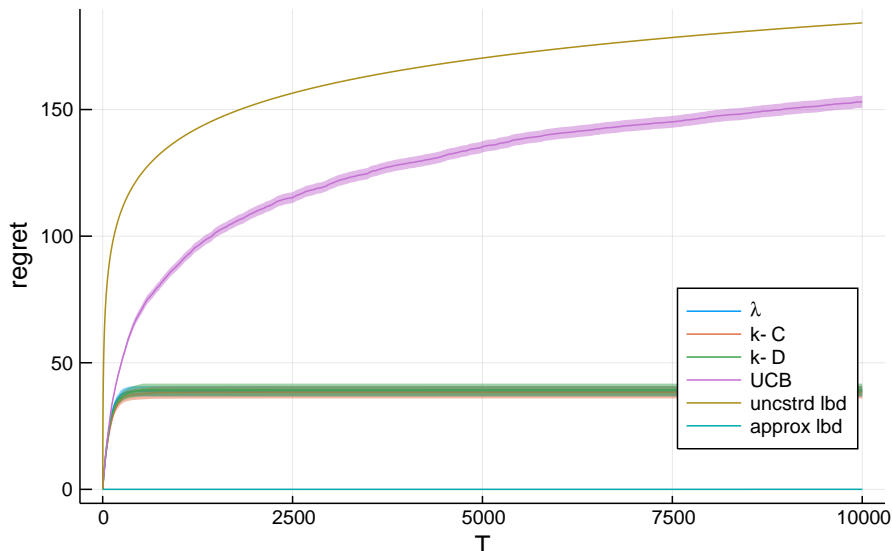
- 1: Pull each arm once and get $\hat{\mu}_K$.
- 2: **for** $t = K + 1, \dots, T$ **do**
- 3: **if** $\exists i \in [K], \min_{\lambda \in -i} \sum_k N_{t-1}^k d(\hat{\mu}_{t-1}^k, \lambda^k) > f(t-1)$ **then**
- 4: $k_t = i$ (if there are several suitable i , pull any one of them)
- 5: **else**
- 6: $\mu_{t-1}^+, j_t = (\arg) \max_{j \in [K]} \hat{\mu}_{t-1}^j + \sqrt{\frac{\ln(n_{t-1}^j, N_{t-1}^j)}{N_{t-1}^j}}$.
- 7: get \tilde{w}_t from learner $\mathcal{A}_{j_t}^k$, compute $w_t^k \propto \tilde{w}_t^k / \tilde{\Delta}^k$.
- 8: compute best response λ_t .
- 9: Compute $\text{UCB}_t^k = \max_{\xi \in [\hat{\mu}_{t-1}^k - \dots, \hat{\mu}_{t-1}^k + \dots]} \left[\frac{d(\xi, \lambda_t^k)}{\max\{\varepsilon_t, \mathbf{1}\{k \neq j_t\}[\mu_{t-1}^+ - \xi]\}} \right]$
- 10: $k_t = \operatorname{argmin}_{k \in [K]} N_{t-1}^k - \sum_{s=1}^t w_s^k$. ▷ Tracking
- 11: **end if**
- 12: Access $X_t^{k_t}$, update $\hat{\mu}_t$ and N_t
- 13: **end for**



Outline

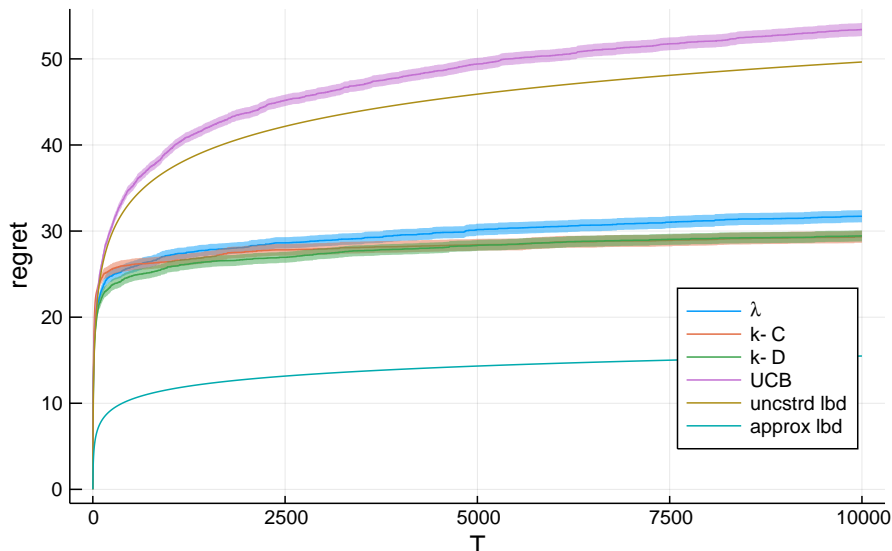
- 1 Introduction
- 2 Lower bound
- 3 Noise Free Case
- 4 The Real Deal
- 5 Experiments**

Experiment: Sparse

 $\mu = [0.3, 0.8, 0.3, 0.3, 0.3, 0.3]$ in Sparse


Experiment: Linear

$\mu = [1.0, 2.21113, 0.366554, -1.98459, -1.5931, 1.0]$ in Linea



Conclusion

Game equilibrium based technique for matching **instance dependent lower bounds** for structured stochastic bandits.

All you need is **Best Response oracle**.

- Fine tuning
- What about “lower-order” terms not scaling with $\ln T$?
- Is minigame interaction “easy data”? MetaGrad [Van Erven and Koolen, 2016]
- Minigames for other problems?

Conclusion

Game equilibrium based technique for matching **instance dependent lower bounds** for structured stochastic bandits.

All you need is **Best Response oracle**.

- Fine tuning
- What about “lower-order” terms not scaling with $\ln T$?
- Is minigame interaction “easy data”? MetaGrad [Van Erven and Koolen, 2016]
- Minigames for other problems?

Thank you!