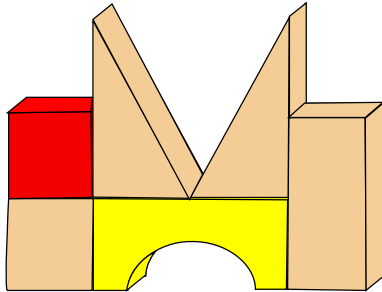


# Sequential Reward Maximisation by Solving a Semi-infinite Covering LP



**Wouter M. Koolen**



Centrum Wiskunde & Informatica

# Team



Rémy Degenne



Han Shao (邵涵)



Wouter Koolen

# Stochastic Bandit



# Stochastic Bandit



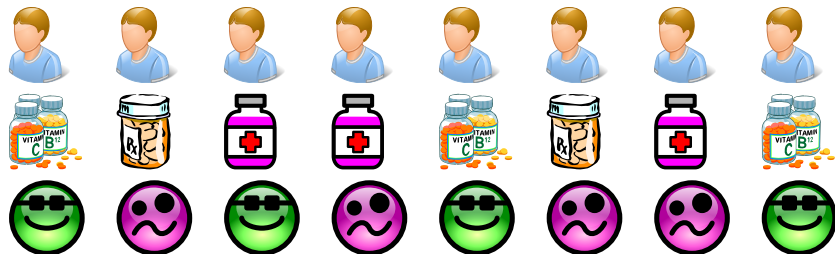
Model (Unknown)

$$\mathbb{P}(\text{Smiley} \mid \text{Rx}) = 1/6$$

$$\mathbb{P}(\text{Smiley} \mid \text{Vitamin}) = 2/3$$

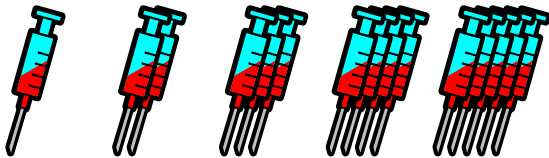
$$\mathbb{P}(\text{Smiley} \mid \text{Cross}) = 1/2$$

# Stochastic Bandit interaction.

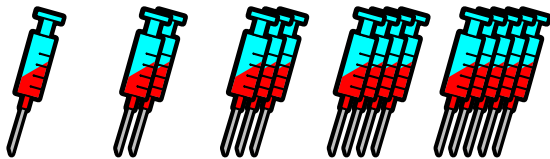


Time →

# Structured Stochastic Bandit



# Structured Stochastic Bandit



## Model (Unknown)

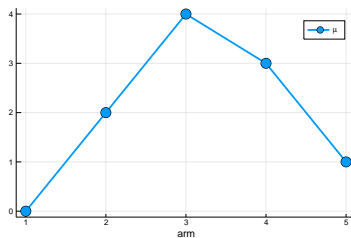
$$\mathbb{P}(\text{smiley} \mid \text{1 syringe} \times 1) = 1/6$$

$$\mathbb{P}(\text{smiley} \mid \text{2 syringes} \times 2) = 3/6$$

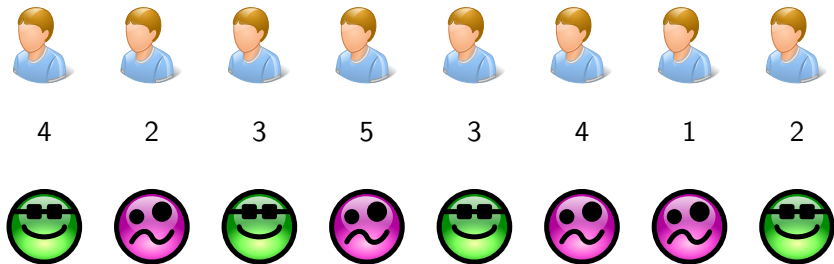
$$\mathbb{P}(\text{smiley} \mid \text{3 syringes} \times 3) = 5/6$$

$$\mathbb{P}(\text{smiley} \mid \text{4 syringes} \times 4) = 4/6$$

$$\mathbb{P}(\text{smiley} \mid \text{5 syringes} \times 5) = 2/6$$

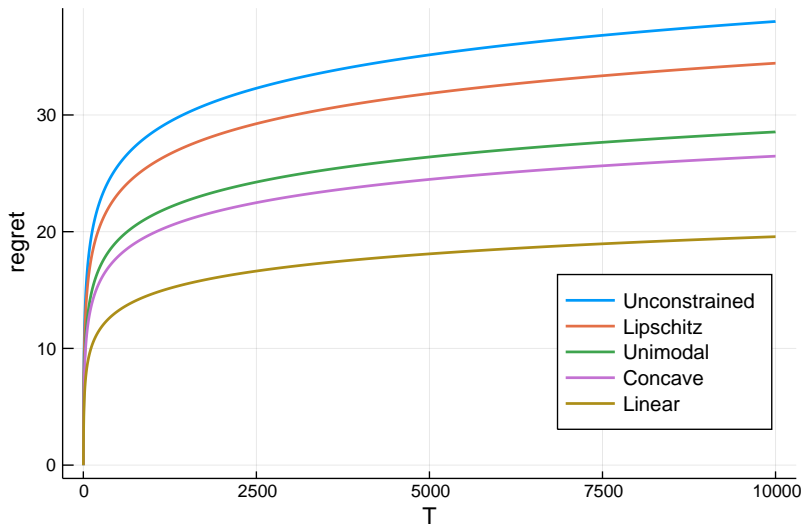
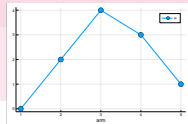


# Structured Stochastic Bandit Interaction





# Desired behaviour

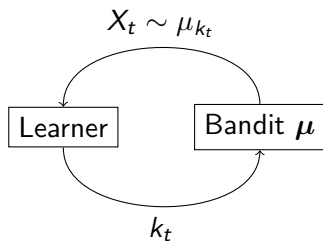




# Outline

- 1 Introduction
- 2 Lower bound
- 3 Noise Free Case
- 4 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

# Bandit 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 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

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

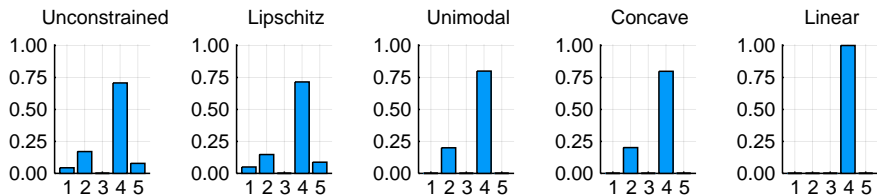
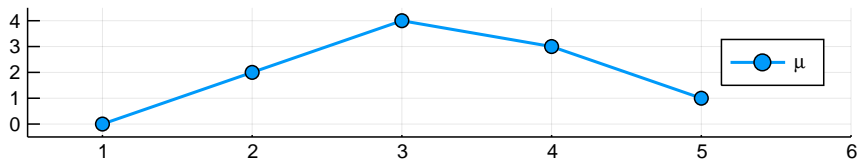
$w^k \propto N^k$   
pulls

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

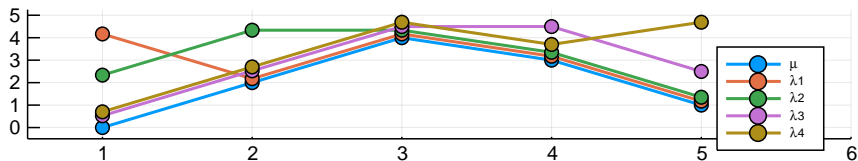
$\tilde{w}^k \propto N^k \Delta^k$   
regret

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

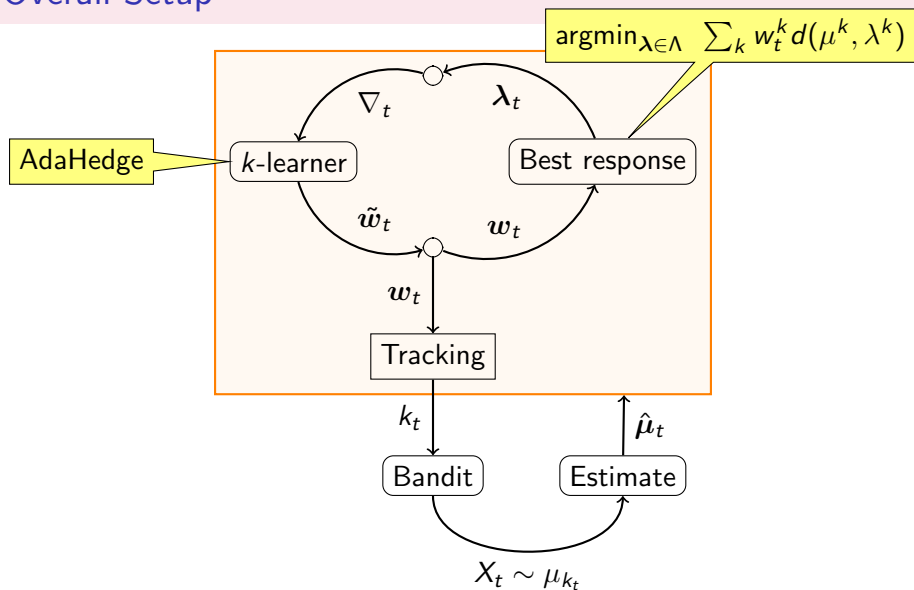
## Illustration



Support for Lipschitz



## Overall Setup



# Outline

- 1 Introduction
- 2 Lower bound
- 3 Noise Free Case**
- 4 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  $\mathcal{B}_n^k = O(\sqrt{n}) = O(\sqrt{\ln T}) = o(\ln T)$ .

# On Symmetry

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

Can also focus on  $\lambda$ -learner instead of  $k$ -learner. Interesting trade-offs

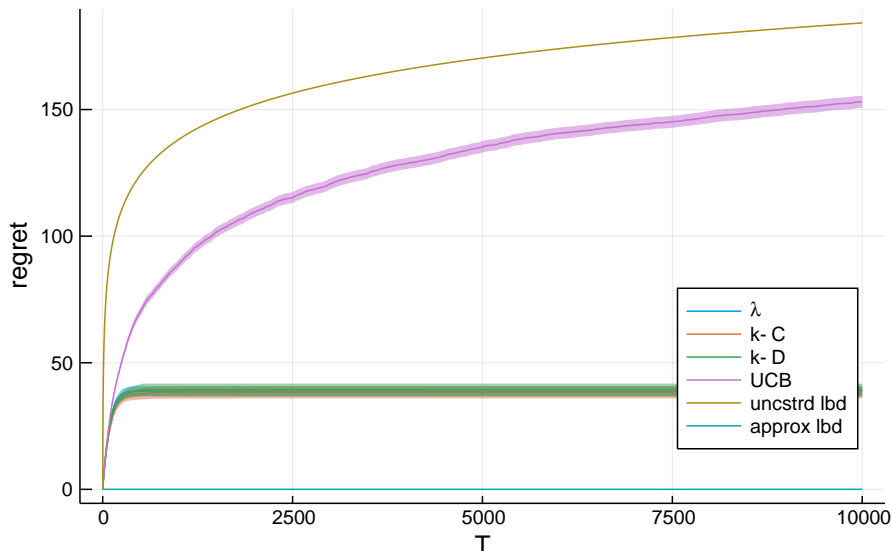
- More complex domain  $\lambda \in \Lambda$ .
- No need for tracking, best response in  $k$  is “pure” arm.

Will show both in experiments.

# Outline

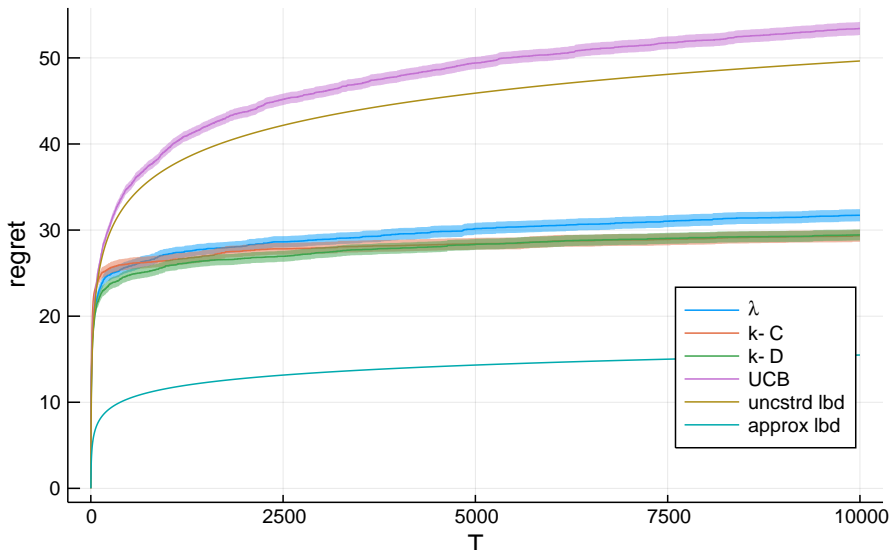
- 1 Introduction
- 2 Lower bound
- 3 Noise Free Case
- 4 Experiments**

## 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!