Exploration and Exploitation in Structured Stochastic Bandits



Wouter M. Koolen



## Collaborators







Rémy Degenne Han Shao (邵涵)

Emilie Kaufmann







# Part I

# Saddle Points Intermezzo





Objective function

g(x, y)

convex in x, concave in y.

### Games



Objective function

g(x, y)

convex in x, concave in y. The game **value** is

$$V^* = \inf_{x} \sup_{y} g(x, y) = \sup_{y} \inf_{x} g(x, y).$$

### Games



Objective function

g(x, y)

convex in x, concave in y. The game **value** is

$$V^* = \inf_{x} \sup_{y} g(x, y) = \sup_{y} \inf_{x} g(x, y).$$

#### Definition

An  $\epsilon$ -saddle point  $(\bar{x}, \bar{y})$  satisfies

$$V^* - \epsilon \leq \inf_{x} g(x, \overline{y}) \leq V^* \leq \sup_{y} g(\overline{x}, y) \leq V^* + \epsilon.$$

Question: how to find  $\epsilon$ -saddle point?

Wouter Koolen



Idea [Freund and Schapire, 1999]: play a regret minimisation algorithm for x against one for y.

- Players play  $x_t$  and  $y_t$ .
- Players see loss functions

$$x \mapsto +g(x, y_t),$$
  
 $y \mapsto -g(x_t, y).$ 

Output pair of **iterate averages**:  $\left(\frac{1}{T}\sum_{t=1}^{T} x_t, \frac{1}{T}\sum_{t=1}^{T} y_t\right)$ .

## Saddle point



Assume the players have regret (bounds)  $R_T^x$  and  $R_T^y$ , i.e.

$$\sum_{t=1}^{T} +g(x_t, y_t) - \inf_{x} \sum_{t=1}^{T} +g(x, y_t) \leq R_T^x$$
$$\sum_{t=1}^{T} -g(x_t, y_t) - \inf_{y} \sum_{t=1}^{T} -g(x_t, y) \leq R_T^y$$

#### Theorem

The iterate averages  $\bar{x}_T = \frac{1}{T} \sum_{t=1}^T x_t$  and  $\bar{y}_T = \frac{1}{T} \sum_{t=1}^T y_t$  form an  $\frac{R_T^x + R_T^y}{T}$ -saddle point.

## Analysis



 $V^* = \inf_{x} \sup_{y} g(x, y)$  $\leq \sup g(\bar{x}_T, y)$ (suboptimal choice  $\bar{x}_T$ )  $\leq \sup_{y} \frac{1}{T} \sum_{t=1}^{T} g(x_t, y)$ (convexity in 1st argument)  $\leq \frac{1}{T}\sum_{t=1}^{T}g(x_t, y_t) + \frac{R_T^y}{T}$ (y player regret guarantee)  $\leq \inf_{x} \frac{1}{T} \sum_{t=1}^{T} g(x, y_t) + \frac{R_T^x + R_T^y}{T}$ (x player regret guarantee)  $\leq \inf_{\mathbf{x}} g(\mathbf{x}, \bar{\mathbf{y}}_T) + \frac{R_T^{\mathbf{x}} + R_T^{\mathbf{y}}}{T}$ (concavity in 2nd argument)  $\leq \inf_{x} \sup_{y} g(x, y) + \frac{R_T^x + R_T^y}{T}$ (suboptimal choice  $\bar{y}_T$ )  $= V^* + \frac{R_T^x + R_T^y}{\tau}$ 







For either player, we can use e.g.

• OGD, Hedge, FTRL, FTPL, ...  $\sqrt{T}$  regret  $\Rightarrow \epsilon = 1/\sqrt{T}$ 



For either player, we can use e.g.

• OGD, Hedge, FTRL, FTPL, ...  $\sqrt{T}$  regret  $\Rightarrow \epsilon = 1/\sqrt{T}$ 

while the second player can ensure negative regret with

- Be-The-Leader
- Best Response no memory



For either player, we can use e.g.

• OGD, Hedge, FTRL, FTPL, ...  $\sqrt{T}$  regret  $\Rightarrow \epsilon = 1/\sqrt{T}$ 

while the second player can ensure negative regret with

- Be-The-Leader
- Best Response no memory

Many more options. Optimism [Rakhlin and Sridharan, 2013], principled path to Nesterov acceleration [Wang and Abernethy, 2018]







## Stochastic Bandit



## Stochastic Bandit



### Model (Unknown)



## Stochastic Bandit Interaction



Time



- O Pure Exploration: use trial to cure population
- Reward Maximisation: cure patients in trial

## Structured Stochastic Bandit



## Structured Stochastic Bandit



#### Model (Unknown)



ntroduction to the Tutoria

## Structured Stochastic Bandit Interaction



## What This Tutorial is About

We will develop **efficient learning algorithms** for **Pure Exploration** and **Reward Maximisation**.

**Information-theoretic lower bounds** will tell us that the complexity of each task is characterised by a certain **two-player zero-sum game**.

We will base our learning algorithms on iterative **saddle point solvers** for this game.

# Part II

# Pure Exploration / Active Testing

## Outline



## 2 Introduction

#### 3 Model

#### 4 Lower Bound

- 5 Pure Exploration Algorithms
- 6 A Games Perspective on TaS
- 7 Experiments

#### 8 Conclusion

## Topic: Pure Exploration





## Topic: Pure Exploration





Main scientific questions

- Efficient systems
- Sample complexity as function of query and environment

## Outline



#### 2 Introduction

#### 3 Model

#### 4 Lower Bound

- 5 Pure Exploration Algorithms
- 6 A Games Perspective on TaS
- 7 Experiments

#### 8 Conclusion

## Environments

We fix an 1-d exponential family (Bernoulli, Gaussian, ...) parameterised by the mean. KL divergence denoted by  $d(\mu, \lambda)$ .

#### Multi-armed bandit model

A K-armed bandit model is a tuple  $\mu = (\mu_1, \dots, \mu_K)$ .

## Environments

We fix an 1-d exponential family (Bernoulli, Gaussian, ...) parameterised by the mean. KL divergence denoted by  $d(\mu, \lambda)$ .

#### Multi-armed bandit model

A K-armed bandit model is a tuple  $\mu = (\mu_1, \dots, \mu_K)$ .

#### Query

Set of possible environments  $\mathcal{M} \subseteq \mathbb{R}^{K}$ . Set of possible answers  $\mathcal{I}$ . **Correct answer** function  $i^* : \mathcal{M} \to \mathcal{I}$ .

#### Mode

### Examples

Problem nameBestPossible answers  $\mathcal{I}$ [K]Correct answer  $i^*(\mu)$ argm

Best Arm [K]argmax<sub>k</sub>  $\mu_k$  Minimum Threshold {lo, hi} lo if min<sub>k</sub>  $\mu_k < \gamma$ hi if min<sub>k</sub>  $\mu_k > \gamma$ 





#### Mode

#### Examples

Problem nameBest ArmPossible answers  $\mathcal{I}$ [K]Correct answer  $i^*(\mu)$  $\operatorname{argmax}_k \mu_k$ 

Minimum Threshold {lo, hi} lo if min<sub>k</sub>  $\mu_k < \gamma$ hi if min<sub>k</sub>  $\mu_k > \gamma$ 





- Top-*M*
- Combinatorial Best Arm
- Maximum Profit
- Unit Ball

- Thresholding Bandit
- Pure Nash equilibrium
- Game Tree Search
- . . .

Wouter Koolen

## Strategy for Learner

#### Strategy

- Stopping rule  $\tau \in \mathbb{N}$
- In round  $t \leq \tau$  sampling rule picks  $A_t \in [K]$ . See  $X_t \sim \mu_{A_t}$ .
- Recommendation rule  $\hat{l} \in \mathcal{I}$ .

## Strategy for Learner

#### Strategy

- Stopping rule  $\tau \in \mathbb{N}$
- In round  $t \leq \tau$  sampling rule picks  $A_t \in [K]$ . See  $X_t \sim \mu_{A_t}$ .
- Recommendation rule  $\hat{l} \in \mathcal{I}$ .

$$\mathsf{Realisation} \,\, \mathsf{of} \,\, \mathsf{interaction} \colon \, \mathcal{H} \coloneqq \Big( \mathsf{A}_1, \mathsf{X}_1, \ldots, \mathsf{A}_\tau, \mathsf{X}_\tau, \widehat{\mathsf{I}} \Big).$$

## Strategy for Learner

#### Strategy

- Stopping rule  $\tau \in \mathbb{N}$
- In round  $t \leq \tau$  sampling rule picks  $A_t \in [K]$ . See  $X_t \sim \mu_{A_t}$ .
- Recommendation rule  $\hat{l} \in \mathcal{I}$ .

Realisation of interaction: 
$$\mathcal{H}\coloneqq \left( \mathsf{A}_1, \mathsf{X}_1, \ldots, \mathsf{A}_{ au}, \mathsf{X}_{ au}, \hat{\mathit{I}} 
ight).$$

**Two** objectives: sample efficiency  $\tau$  and correctness  $\hat{l} = i^*(\mu)$ .

## Goal: PAC learning



#### Definition

Fix small confidence  $\delta \in (0, 1)$ . A strategy is  $\delta$ -correct if

$$\mathbb{P}_{oldsymbol{\mu}}ig(\hat{l}
eq i^*(oldsymbol{\mu})ig)\ \leq\ \delta$$
 for every bandit model  $oldsymbol{\mu}\in\mathcal{M}.$ 

## Goal: PAC learning



#### Definition

Fix small confidence  $\delta \in (0, 1)$ . A strategy is  $\delta$ -correct if

$$\mathbb{P}_{oldsymbol{\mu}}ig(\hat{m{l}}
eq i^*(oldsymbol{\mu})ig) \ \leq \ \delta$$
 for every bandit model  $oldsymbol{\mu}\in\mathcal{M}.$ 

Goal: minimise sample complexity  $\mathbb{E}_{\mu}[\tau]$  over all  $\delta$ -correct strategies.
# Goal: PAC learning



#### Definition

Fix small confidence  $\delta \in (0, 1)$ . A strategy is  $\delta$ -correct if

$$\mathbb{P}_{oldsymbol{\mu}}ig(\hat{m{l}}
eq i^*(oldsymbol{\mu})ig) \ \leq \ \delta$$
 for every bandit model  $oldsymbol{\mu}\in\mathcal{M}.$ 

Goal: minimise sample complexity  $\mathbb{E}_{\mu}[\tau]$  over all  $\delta$ -correct strategies.

Not in this talk: Fixed Budget.

# Outline



### 2 Introduction

#### 3 Model

### 4 Lower Bound

- 5 Pure Exploration Algorithms
- 6 A Games Perspective on TaS
- 7 Experiments

#### 8 Conclusion

### Instance-Dependent Sample Complexity Lower Bound

Intuition (going back at least to Lai and Robbins [1985]): if observations are likely under both  $\mu$  and  $\lambda$ , yet  $i^*(\mu) \neq i^*(\lambda)$ , then learner cannot stop and be correct in both.

### Instance-Dependent Sample Complexity Lower Bound

Intuition (going back at least to Lai and Robbins [1985]): if observations are likely under both  $\mu$  and  $\lambda$ , yet  $i^*(\mu) \neq i^*(\lambda)$ , then learner cannot stop and be correct in both.

Define the **alternative** to answer  $i \in \mathcal{I}$  by  $\neg i \coloneqq \{\lambda \in \mathcal{M} | i^*(\lambda) \neq i\}$ .

### Instance-Dependent Sample Complexity Lower Bound

Intuition (going back at least to Lai and Robbins [1985]): if observations are likely under both  $\mu$  and  $\lambda$ , yet  $i^*(\mu) \neq i^*(\lambda)$ , then learner cannot stop and be correct in both.

Define the **alternative** to answer  $i \in \mathcal{I}$  by  $\neg i \coloneqq \{\lambda \in \mathcal{M} | i^*(\lambda) \neq i\}$ .

Theorem (Castro 2014, Garivier and Kaufmann 2016) Fix a  $\delta$ -correct strategy. Then for every bandit model  $\mu \in \mathcal{M}$ 

$$\mathbb{E}_{oldsymbol{\mu}}[ au] \ \geq \ {\mathcal T}^*(oldsymbol{\mu}) \, \ln rac{1}{\delta}$$

where the characteristic time  $T^*(\mu)$  is given by

$$rac{1}{\mathcal{T}^*(oldsymbol{\mu})} = \max_{oldsymbol{w}\in riangle_K} \min_{oldsymbol{\lambda}\in 
eg i^*(oldsymbol{\mu})} \sum_{i=1}^K w_i d(\mu_i,\lambda_i).$$

### Example

Best Arm identification:  $i^*(\mu) = \operatorname{argmax}_i \mu_i$ . K = 5 Bernoulli arms,  $\mu = (0.4, 0.3, 0.2, 0.1, 0.0)$ .

$$T^*(\mu) = 200.4$$
  $w^*(\mu) = (0.45, 0.46, 0.06, 0.02, 0.01)$ 

At  $\delta = 0.05$ , the time gets multiplied by  $\ln \frac{1}{\delta} = 3.0$ .

# Outline



### 2 Introduction

#### 3 Model

#### 4 Lower Bound

5 Pure Exploration Algorithms

6 A Games Perspective on TaS

#### 7 Experiments

#### 8 Conclusion

Pure Exploration Algorithms

Recall, a strategy is defined by

- Stopping rule
- Recommendation rule
- Sampling rule

# Stopping and Recommendation Rules

Quantify evidence for  $\{i^*(\mu) = i^*(\hat{\mu}_t)\}$  vs alternative  $\{i^*(\mu) \neq i^*(\hat{\mu}_t)\}$ .

# Stopping and Recommendation Rules

Quantify evidence for  $\{i^*(\mu) = i^*(\hat{\mu}_t)\}$  vs alternative  $\{i^*(\mu) \neq i^*(\hat{\mu}_t)\}$ . Definition

The extended GLR statistic is defined as

$$\hat{\Lambda}_t = \inf_{\lambda \in \neg i^*(\hat{\mu}(t))} \sum_{a=1}^K N_a(t) d\left(\hat{\mu}_a(t), \lambda_a\right).$$

# Stopping and Recommendation Rules

Quantify evidence for  $\{i^*(\mu) = i^*(\hat{\mu}_t)\}$  vs alternative  $\{i^*(\mu) \neq i^*(\hat{\mu}_t)\}$ . Definition

The extended GLR statistic is defined as

$$\hat{\Lambda}_t = \inf_{\lambda \in \neg i^*(\hat{\mu}(t))} \sum_{a=1}^{K} N_a(t) d\left(\hat{\mu}_a(t), \lambda_a\right).$$

Proposal: stop when

$$\hat{\Lambda}_t \geq \beta(t, \delta) \coloneqq \ln \frac{1}{\delta} + K \ln \ln t + K \ln \ln \frac{1}{\delta}$$
  
and recommend  $\hat{l} = i^*(\hat{\mu}_t)$ .

### Theorem (Kaufmann and Koolen 2018)

The above stopping and recommendation rules, combined with any sampling rule give a  $\delta$ -correct algorithm.

Wouter Kooler

Bandits, Games, Explore/Exploit

WLT2 31/83

# Operationalisation of the Oracle Weights

Recall sample complexity lower bound governed by

$$\max_{\boldsymbol{w} \in \triangle_{K}} \min_{\boldsymbol{\lambda} \in \neg i^{*}(\boldsymbol{\mu})} \sum_{i=1}^{K} w_{i} d(\mu_{i}, \lambda_{i})$$

# Operationalisation of the Oracle Weights

Recall sample complexity lower bound governed by

$$\max_{\boldsymbol{w} \in \triangle_{\kappa}} \min_{\boldsymbol{\lambda} \in \neg i^{*}(\boldsymbol{\mu})} \sum_{i=1}^{\kappa} w_{i} d(\mu_{i}, \lambda_{i})$$

Any matching algorithm must sample with optimal (oracle) proportions

$$m{w}^*(m{\mu}) = rgmax_{m{w}\in riangle_K} \min_{m{\lambda}\in 
eg i^*(m{\mu})} \sum_{i=1}^K w_i d(\mu_i, \lambda_i)$$

Idea: draw  $A_t \sim w^*(\hat{\mu}(t))$ .

Idea: draw  $A_t \sim w^*(\hat{\mu}(t)).$ 

Track-and-Stop [Garivier and Kaufmann, 2016]

- Ensure  $\hat{\mu}(t) 
  ightarrow \mu$  by forced exploration
- ullet assuming  $w^*$  is continuous, this ensures  $w^*(\hat{\mu}_t) o w^*(\mu)$ .
- Draw arm with  $N_i(t)$  below  $\sum_{s=1}^t w_i^*(\hat{\mu}_s)$  (C-tracking)

• hence 
$$N_i(t)/t o w_i^*(oldsymbol{\mu})$$

Inherit  $\delta$ -correctness from GLR stopping/recommendation rule.

Idea: draw  $A_t \sim w^*(\hat{\mu}(t)).$ 

Track-and-Stop [Garivier and Kaufmann, 2016]

- Ensure  $\hat{\mu}(t) 
  ightarrow \mu$  by forced exploration
- ullet assuming  $w^*$  is continuous, this ensures  $w^*(\hat{\mu}_t) o w^*(\mu)$ .
- Draw arm with  $N_i(t)$  below  $\sum_{s=1}^t w_i^*(\hat{\mu}_s)$  (C-tracking)
- hence  $N_i(t)/t o w_i^*(\mu)$

Inherit  $\delta$ -correctness from GLR stopping/recommendation rule.

### Theorem (Degenne et al. 2019)

*Track-and-Stop with C-tracking has asymptotically optimal sample complexity.* 

Idea: draw  $A_t \sim w^*(\hat{\mu}(t)).$ 

Track-and-Stop [Garivier and Kaufmann, 2016]

- $\bullet$  Ensure  $\hat{\mu}(t) \rightarrow \mu$  by forced exploration
- ullet assuming  $w^*$  is continuous, this ensures  $w^*(\hat{\mu}_t) o w^*(\mu)$ .
- Draw arm with  $N_i(t)$  below  $\sum_{s=1}^t w_i^*(\hat{\mu}_s)$  (C-tracking)
- hence  $N_i(t)/t 
  ightarrow w_i^*(oldsymbol{\mu})$

Inherit  $\delta$ -correctness from GLR stopping/recommendation rule.

### Theorem (Degenne et al. 2019)

*Track-and-Stop with C-tracking has asymptotically optimal sample complexity.* 

### Theorem (Degenne et al. 2019)

Track-and-Stop with D-tracking may fail to converge.

# Outline



### 2 Introduction

#### 3 Model

#### 4 Lower Bound

#### 5 Pure Exploration Algorithms

#### 6 A Games Perspective on TaS

#### 7 Experiments

#### 8 Conclusion

### Games perspective

Recall TaS based on plug-in estimate  $w^*(\hat{\mu}_t)$  of oracle weights

$$m{w}^*(m{\mu}) = rgmax \min_{m{w}\in riangle_K} \sum_{m{\lambda}\in 
eg i^*(m{\mu})}^K w_i d(\mu_i, \lambda_i)$$

### Games perspective

Recall TaS based on plug-in estimate  $w^*(\hat{\mu}_t)$  of oracle weights

$$m{w}^*(m{\mu}) \;=\; rgmax_{m{w}\in riangle_K} \min_{m{\lambda}\in 
eg i^*(m{\mu})} \; \sum_{i=1}^K w_i d(\mu_i,\lambda_i)$$

We can implement the Track-and-Stop sampling rule by running an online-learning based saddle point solver to (approximate) convergence every round.

Choice of learners: AdaHedge vs Best Response.

### Games perspective

Recall TaS based on plug-in estimate  $w^*(\hat{\mu}_t)$  of oracle weights

$$m{w}^*(m{\mu}) \;=\; rgmax_{m{w}\in riangle_K} \min_{m{\lambda}\in 
eg i^*(m{\mu})} \; \sum_{i=1}^K w_i d(\mu_i,\lambda_i)$$

We can implement the Track-and-Stop sampling rule by running an online-learning based saddle point solver to (approximate) convergence every round.

Choice of learners: AdaHedge vs Best Response.

The user needs to provide **best response oracle** (often tractable)

$$w, \mu \mapsto \operatorname*{argmin}_{\lambda \in \neg i^*(\mu)} \sum_{i=1}^{K} w_i d(\mu_i, \lambda_i).$$

## Ironing out Inefficiencies

- We are computing  $w^*(\hat{\mu}_t)$  on noisy  $\hat{\mu}_t \approx \mu$ . So how precise does our saddle point need to be?
- $\hat{\mu}_t pprox \hat{\mu}_{t+1}$ . Can we reuse (most) computation?
- Can we get finite confidence guarantees?

Interleaved Iterative Solution

Main idea [Degenne, Koolen, and Ménard, 2019]: advance the saddle point solver by **one** iteration for every bandit interaction.

# Sampling Rule



# Compositionality

- The "overheads" of the ingredients **compose**: Tracking O(1), concentration  $\sqrt{T}$ , regret  $\sqrt{T}$ , optimism  $\sqrt{T}$ .
- Theorem (Degenne, Koolen, and Ménard 2019) The sample complexity is at most

$$\mathbb{E}_{oldsymbol{\mu}}[ au] \ \le \ extsf{T}^*(oldsymbol{\mu}) \ {\sf ln} \ rac{1}{\delta} + {\it small}$$

### Proof ideas (cheating with optimism, $i_t = i^*$ )

As long as we do not stop,  $t < \tau$ ,

$$\begin{split} \beta(t,\delta) &\geq \inf_{\boldsymbol{\lambda} \in \neg i_t} \sum_{k=1}^{K} N_t^k d(\mu^k, \lambda^k) \qquad (\text{stop rule}) \\ &\approx \inf_{\boldsymbol{\lambda} \in \neg i^*} \sum_{s=1}^t \sum_{k=1}^K w_s^k d(\mu^k, \lambda^k) \qquad (\text{tracking}) \\ &\geq \sum_{s=1}^t \sum_{k=1}^K w_s^k \mathbb{E}_{\boldsymbol{\lambda} \sim q_s} d(\mu^k, \lambda^k) - R_t^{\boldsymbol{\lambda}} \qquad (\text{regret } \boldsymbol{\lambda}) \\ &\geq \max_k \sum_{s=1}^t \mathbb{E}_{\boldsymbol{\lambda} \sim q_s} d(\mu^k, \lambda^k) - R_t^{\boldsymbol{\lambda}} - R_t^k \qquad (\text{regret } k) \\ &\geq t \inf_{q \in \mathcal{P}(\neg i^*)} \max_k \mathbb{E}_{\boldsymbol{\lambda} \sim q} d(\mu^k, \lambda^k) - O(\sqrt{t}) \end{split}$$

Find maximal t to get bound on  $\tau$ .

Wouter Koole

# Outline



### 2 Introduction

#### 3 Model

#### 4 Lower Bound

- 5 Pure Exploration Algorithms
- 6 A Games Perspective on TaS

#### 7 Experiments

#### **Conclusion**

### Best Arm Identification Experiment

Best Arm for Bernoulli bandit model  $\mu = (0.3, 0.21, 0.2, 0.19, 0.18)$ . The oracle weights are  $w^* = (0.34, 0.25, 0.18, 0.13, 0.10)$ .  $\delta = 0.1$ .



### Minimum Threshold Experiment

Minimum Threshold for Gaussian bandit model  $\mu = (0.5, 0.6)$  with threshold  $\gamma = 0.6$ ,  $w^* = (1, 0)$ . Note the excessive sample complexity of T-C/T-D.  $\delta = 10^{-10}$ .



# Outline



### 2 Introduction

- 3 Model
- 4 Lower Bound
- 5 Pure Exploration Algorithms
- 6 A Games Perspective on TaS
- 7 Experiments



# Conclusion

- Pure Exploration currently going through a renaissance
- New and different instance-optimal identification algorithms
  - Best Arm
  - Combinatorial best action
  - Game Tree Search
  - ▶ ...
- Moving toward more complex queries. RL on the horizon ....
- Useful submodules

# **Topics Skipped**

- Multiple correct answers (e.g. *ϵ* Best Arm) [Degenne and Koolen, 2019] ⇐ surprisingly subtle.
- Optimal algorithms based on variations of Thompson Sampling
  - Top-Two for Best Arm [Russo, 2016]
  - Murphy Sampling for Minimum Threshold [Kaufmann et al., 2018].

# Many questions remain open

- Practically efficient algorithms
- Remove forced exploration
- Moderate confidence  $\delta \not\rightarrow 0$  regime [Simchowitz et al., 2017].
- Understand sparsity patterns
- Dynamically expanding horizon

# Part III

# Reward Maximisation / Regret Minimisation

Now samples represent **reward**, and the algorithm aims to collect **as much reward as possible**.

Regret metric: maximum reward minus reward collected.

Famous UCB algorithm (family) [Auer et al., 2002].

 $O(\ln T)$  regret possible.

# Stochastic Bandit Instance (Running Example)


## Desired behaviour





### Outline





- 10 Lower bound
- 11 Noise Free Case
- 12 The Real Deal



Introduction

### Setting





Structure  $\mathcal{M} \subseteq R^{\mathcal{K}}$ . MAB instance  $\mu \in \mathcal{M}$ Time horizon TExpfam  $d(\mu, \lambda)$ Gaps  $\Delta^k = \mu^* - \mu^k$ 

Regret := 
$$\sum_{k=1}^{K} \mathbb{E}[N_T^k] \Delta^k$$

#### 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] part 1

### Outline



#### Introduction



Noise Free Case





### 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  $i^*(\mu) \neq i^*(\lambda)$ :

$$\mathbb{E}_{\mu}[N_{T}^{i^{*}(\mu)}]/T \to 1$$
 but  $\mathbb{E}_{\lambda}[N_{T}^{i^{*}(\mu)}]/T \to 0.$ 

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

$$\mathsf{KL}(\mathbb{P}^{\mathcal{T}}_{\boldsymbol{\mu}} \, \big\| \, \mathbb{P}^{\mathcal{T}}_{\boldsymbol{\lambda}}) = \sum_{k} \mathbb{E}_{\boldsymbol{\mu}}[N^{k}_{\mathcal{T}}] d(\mu^{k}, \lambda^{k}) \geq \ln \mathcal{T}.$$

### Instance-Dependent Regret Lower Bound



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

$$V_{\mathcal{T}} = \min_{N \ge 0} \sum_{k} N^{k} \Delta^{k} \text{ subject to } \inf_{\lambda \in \Lambda} \sum_{k} N^{k} d(\mu^{k}, \lambda^{k}) \ge \ln \mathcal{T}$$

where

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

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

### Operationalising the Lower Bound

Earlier work

At each time step

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

## Operationalising the Lower Bound

Earlier work

At each time step

- ullet compute oracle sample counts  $N^*(\hat{\mu}_t)$  and advance  $N_t o N^*$ , 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

$$D^* = \underbrace{\max_{w \in \Delta} \inf_{\lambda \in \Lambda} \frac{\sum_k w^k d(\mu^k, \lambda^k)}{\sum_k w^k \Delta^k}}_{\substack{w^k \propto N^k \\ \text{pulls}}}$$

....

# Minigame

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

$$D^{*} = \underbrace{\max_{w \in \Delta} \inf_{\lambda \in \Lambda} \frac{\sum_{k} w^{k} d(\mu^{k}, \lambda^{k})}{\sum_{k} w^{k} \Delta^{k}}}_{\substack{w^{k} \propto N^{k}}}$$
$$= \underbrace{\max_{w \in \Delta} \inf_{\lambda \in \Lambda} \sum_{k} \tilde{w}^{k} \frac{d(\mu^{k}, \lambda^{k})}{\Delta^{k}}}_{\substack{w^{k} \propto N^{k} \Delta^{k}}}$$

k

NIL

# Minigame

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

$$D^{*} = \underbrace{\max_{w \in \Delta} \inf_{\lambda \in \Lambda} \frac{\sum_{k} w^{k} d(\mu^{k}, \lambda^{k})}{\sum_{k} w^{k} \Delta^{k}}}_{\text{pulls}}$$
$$= \underbrace{\max_{\tilde{w} \in \Delta} \inf_{\lambda \in \Lambda} \sum_{k} \tilde{w}^{k} \frac{d(\mu^{k}, \lambda^{k})}{\Delta^{k}}}_{\tilde{w} \in \Delta} \underbrace{\tilde{w}^{k} \propto N^{k} \Delta^{k}}_{\text{regret}}$$
$$= \underbrace{\inf_{q \in \Delta(\Lambda)} \max_{k} \frac{\mathbb{E}_{\lambda \sim q} \left[ d(\mu^{k}, \lambda^{k}) \right]}{\Delta^{k}}}_{\Delta^{k}}$$

- - 1

Lower bound

#### Illustration



#### **Overall Setup**



### Outline













### 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  $w_{1}, \ldots, w_{n}$  satisfy

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

Note

- Can get  $k_1, \ldots, k_n$  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  $w_t \in riangle_K$  and  $\lambda_t \in \Lambda$ , we instantiate a *k*-learner for the gain function

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

to provide regret bound

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

# Regret analysis (ctd)

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

$$w_t^k \propto \tilde{w}_t^k / \Delta^k$$
 (2)  
 $\lambda_t \in \operatorname{argmin}_{\lambda \in \Lambda} \sum_k w_t^k d(\mu^k, \lambda^k)$  (3)

to obtain

$$\sum_{t=1}^{n} g_{t}(\tilde{w}_{t}) = \sum_{t=1}^{n} \langle w_{t}, \Delta \rangle \sum_{k} \tilde{w}_{t}^{k} \frac{d(\mu^{k}, \lambda_{t}^{k})}{\Delta^{k}} \stackrel{(2)}{=} \sum_{t=1}^{n} \sum_{k} w_{t}^{k} d(\mu^{k}, \lambda_{t}^{k})$$
$$\stackrel{(3)}{=} \sum_{t=1}^{n} \inf_{\lambda \in \Lambda} \sum_{k} w_{t}^{k} d(\mu^{k}, \lambda^{k}) \leq \inf_{\lambda \in \Lambda} \sum_{t=1}^{n} \sum_{k} w_{t}^{k} d(\mu^{k}, \lambda^{k})$$
(4)



# Regret analysis (ctd)



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

$$\ln T + \mathcal{B}_{n}^{k} \geq \max_{k} \sum_{t=1}^{n} \langle w_{t}, \Delta \rangle \frac{d(\mu^{k}, \lambda_{t}^{k})}{\Delta^{k}} = R_{n} \max_{k} \sum_{t=1}^{n} \frac{\langle w_{t}, \Delta \rangle}{R_{n}} \frac{d(\mu^{k}, \lambda_{t}^{k})}{\Delta^{k}}$$
$$\geq R_{n} \inf_{q \in \Delta(\Lambda)} \max_{k} \frac{\mathbb{E}_{\lambda \sim q} \left[ d(\mu^{k}, \lambda^{k}) \right]}{\Delta^{k}} = R_{n} D^{*}$$

where we abbreviated  $R_n = \sum_{t=1}^n \langle w_t, \Delta \rangle$ . 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  $\lambda$ -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



#### Introduction

- 10 Lower bound
- Noise Free Case



#### 13 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  $\mu$ .
- 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

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

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  $\mu$  by estimate  $\hat{\mu}_t$ .
- Add optimism to force exploration.
   We introduce upper confidence bounds on the ratio KL/gap.

$$\begin{aligned} \text{UCB}_{s}^{k} &= \sup_{\xi \in \mathcal{C}_{s-1}^{k}} \frac{d(\xi, \lambda_{t}^{k})}{\max\left\{\epsilon_{s}, \mathbf{1}\{k \neq j_{s}\}\left[\mu_{s-1}^{+} - \xi\right]\right\}} \end{aligned}$$
  
where  $\mathcal{C}_{s-1}^{k} &= \left[\hat{\mu}_{s-1}^{k} \pm \sqrt{\frac{\overline{\ln}(n_{s-1}^{j_{s}}, N_{s-1}^{k})}{N_{s-1}^{k}}}\right]. \end{aligned}$ 

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  $\exists i \in [K]$ ,  $\min_{\lambda \in \neg i} \sum_{k} N_{t-1}^{k} d(\hat{\mu}_{t-1}^{k}, \lambda^{k}) > f(t-1)$  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  $\lambda_{t}$ .  
9: Compute 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\_{\xi \in t, 1\{k \neq j\_{t}\}[\mu\_{t-1}^{k} - \xi]\}} \right]  
10:  $k_{t} = \arg\min_{k \in [K]} N_{t-1}^{k} - \sum_{s=1}^{t} w_{s}^{k}$ .  $\triangleright$  Tracking  
11: end if  
12: Access  $X_{t}^{k_{t}}$ , update  $\hat{\mu}_{t}$  and  $N_{t}$ 

### Outline



#### Introduction

- 10 Lower bound
- 11 Noise Free Case





#### Experiment: Sparse



#### Experiment: Linear



#### 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 In 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 In T?
- Is minigame interaction "easy data"? MetaGrad [van Erven and Koolen, 2016]
- Minigames for other problems?

# Thank you!





#### Discontinuous single-answer problems

#### 15 KL Tensorises

Wouter Koolen

## About that continuity assumption?

Can  $w^*$  be discontinuous?
# About that continuity assumption?

Can  $w^*$  be discontinuous?

Example: Minimum Threshold



## Continuity restored

Recall oracle weights are given by

$$m{w}^*(m{\mu}) = rgmax_{m{w}\in riangle} \inf_{m{\lambda}\in 
eg i^*(m{\mu})} \sum_{m{a}} w_{m{a}} d(\mu_{m{a}},\lambda_{m{a}})$$

# Continuity restored

Recall oracle weights are given by

$$w^*(\mu) = rgmax_{oldsymbol{w}\in riangle} \inf_{oldsymbol{\lambda}\in 
eg i^*(\mu)} \sum_{oldsymbol{a}} w_{oldsymbol{a}} d(\mu_{oldsymbol{a}},\lambda_{oldsymbol{a}})$$

#### Theorem

 $w^*$ , when viewed as a **set-valued** function, is upper hemicontinuous. Moreover, its output is always a convex set.

### Intuition



On bandit model  $\mu_{\text{i}}$  our empirical distribution will be a convex combination of  $\delta_1$  and  $\delta_2.$ 

## Outline



#### Discontinuous single-answer problems



# **KL** Tensorises

Algorithm is common and observations are IID. Hence

$$\frac{P_{\mu}(I^{T}, X^{T})}{P_{\lambda}(I^{T}, X^{T})} = \prod_{t=1}^{T} \frac{P_{Alg}(I_{t}|I^{t-1}, X^{t-1}) \mu_{l_{t}}(X_{t})}{P_{Alg}(I_{t}|I^{t-1}, X^{t-1}) \lambda_{l_{t}}(X_{t})} = \prod_{t=1}^{T} \frac{\mu_{l_{t}}(X_{t})}{\lambda_{l_{t}}(X_{t})}$$

It follows that

$$\begin{aligned} \mathsf{KL}\big(P_{\mu}(I^{T}, X^{T})\big\|P_{\lambda}(I^{T}, X^{T})\big) &= \sum_{t=1}^{T} \mathbb{E}_{\mu}\left[\ln\frac{\mu_{I_{t}}(X_{t})}{\lambda_{I_{t}}(X_{t})}\right] \\ &= \sum_{t=1}^{T} \mathbb{E}_{\mu}\left[d(\mu_{I_{t}}, \lambda_{I_{t}})\right] \\ &= \sum_{i=1}^{K} \mathbb{E}_{\mu}\left[N_{i,T}\right]d(\mu_{i}, \lambda_{i}) \end{aligned}$$