## **Pure Exploration Problems**

Wouter M. Koolen CWI and University of Twente

MML, June 20, 2024

## Part I

# **Motivating Problems**

I'll show you three motivating applications:

- A/B Testing
- Self-driving
- Solving Games (MCTS)

I'll show you three motivating applications:

- A/B Testing
- Self-driving
- Solving Games (MCTS)

Here is what is common:

- The learner needs to identify some aspect of reality
- The learner can choose which data to collect

## A/B Testing





- A/B testing
- Adaptive clinical trial
- Best arm identification

#### Let's find out in production!





## **Self-Driving**

We are excited about autonomous driving.

We are excited about autonomous driving.

Training AI systems (e.g. deep neural networks) takes a lot of data.

We are excited about autonomous driving.

Training AI systems (e.g. deep neural networks) takes a lot of data.

Yet not all data are equally valuable/useful.

We are excited about autonomous driving.

Training AI systems (e.g. deep neural networks) takes a lot of data.

Yet not all data are equally valuable/useful.

Let's optimise and automate the data collection.

We are excited about autonomous driving.

Training AI systems (e.g. deep neural networks) takes a lot of data.

Yet not all data are equally valuable/useful.

Let's optimise and automate the data collection.

Where should I send my prototype for training?) to find the best driving strategy?





Where should I send my prototype for training to find the best driving strategy ?

- Al system is huge parameterised model
- Lots of possible environments to drive in
- Multiple objectives (safety, efficiency, ...)
- Feedback (crash/intervention) is very one-sided



Where should I send my prototype for training to find the best driving strategy?

- Al system is huge parameterised model
- Lots of possible environments to drive in
- Multiple objectives (safety, efficiency, ...)
- Feedback (crash/intervention) is very one-sided

#### Distilled goal:

• Identify Overall safest: fewest crashes in natural environment mix.

Close-by parameters in close-by environments result in close-by outcomes

Close-by parameters in close-by environments result in close-by outcomes

Close-by parameters in close-by environments result in close-by outcomes



Close-by parameters in close-by environments result in close-by outcomes



Close-by parameters in close-by environments result in close-by outcomes



The world simplifies to vector + table:

The world simplifies to vector + table:

Known natural environment mix



The world simplifies to vector + table:

Known natural environment mix



Unknown crash probabilities



The world simplifies to vector + table:

Known natural environment mix

Unknown crash probabilities



Together these determine the best parameter on average. Say *compared*.

































## Solving Games (MCTS)

**BAI** again



Unknown environment



How to find the best arm by sampling arms?

### Game DAGs




# And Many More

# Beyond best arm

#### Many possible desiderata

- Find the best *M* arms
- Find all sufficiently good arms
- Arms with combinatorial structure
- Continuously many arms (bandit optimisation)
- Prior knowledge about structure (dosage)
- Personalisation (find a policy: context  $\rightarrow$  arm)
- Multi-objective problems
- Risk measures
- Customers return (!) (MDP)

#### and challenges

- delays
- censoring
- constraints
- non-stationarity
- comparison feedback

# Booming Industry within Multi-Armed Bandit / Testing Literature



#### Explosion of bandit testing papers with analogous problems:

- Top-*m* 2017; 2017
- Spectral 2021
- Stratified 2021
- Lipschitz 2019
- Linear 2020; 2020

- Threshold 2017
- MaxGap 2019
- Duelling 2021
- Contextual 2020; 2020
- Pareto 2023

- Minimum 2018
- MCTS 2016
- Markov 2019
- Tail-Risk 2021
- MDP 2021

# Part II

# **Theory and Algorithms**

# Setup



K-armed bandit  $\mu = (\mu_1, \ldots, \mu_K)$ .



K-armed bandit  $\mu = (\mu_1, \ldots, \mu_K)$ .

Each arm k represented by a Bernoulli rate  $\mu_k$  in [0,1].



K-armed bandit  $\boldsymbol{\mu} = (\mu_1, \ldots, \mu_K).$ 

# Each arm k represented by a Bernoulli rate $\mu_k$ in [0, 1].

Observations from arm k are i.i.d. Bernoulli $(\mu_k)$ ).



K-armed bandit  $\boldsymbol{\mu} = (\mu_1, \ldots, \mu_K).$ 

Each arm k represented by a Bernoulli rate  $\mu_k$  in [0, 1].

Observations from arm k are i.i.d. Bernoulli $(\mu_k)$ ).

The best action at the root is:

 $i^*(\mu) \coloneqq \arg \max_i \min_j \max_k \mu_{\mathsf{leaf}(i,j,k)}$ 

# Protocol

We work in the setting of fixed confidence  $\delta \in (0, 1)$ .

# Protocol For $t = 1, 2, ..., \tau$ : • Learner picks an arm $I_t \in [K]$ . • Learner sees $X_t \sim Bernoulli(\mu_{I_t})$ Learner recommends child of the root $\hat{\imath}$

# Protocol

We work in the setting of fixed confidence  $\delta \in (0, 1)$ .

# Protocol For $t = 1, 2, ..., \tau$ : • Learner picks an arm $l_t \in [K]$ . • Learner sees $X_t \sim Bernoulli(\mu_{l_t})$ Learner recommends child of the root $\hat{\imath}$

Strategy for Learner specified by

- sampling rule  $I_t$
- stopping rule  $\tau$
- recommendation rule  $\hat{\imath}$



#### Learner is $\delta$ -correct if for any bandit instance $\mu$

$$\mathbb{P}_{oldsymbol{\mu}}\left\{ au < \infty \wedge \hat{\imath} 
eq i^{*}(oldsymbol{\mu})
ight\} \;\leq\; \delta$$

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

# Sample Complexity Lower Bound



Suppose we believe we are in bandit  $\mu$ .

But we are required to be  $\delta$  correct in all bandits.

When can we stop and recommend answer  $i \equiv i^*(\mu)$ ?



Suppose we believe we are in bandit  $\mu$ .

But we are required to be  $\delta$  correct in all bandits.

When can we stop and recommend answer  $i \equiv i^*(\mu)$ ?

Hypothesis test! Composite null vs point alternative:

$$\mathcal{H}_0 = \neg i \coloneqq ig\{ oldsymbol{\lambda} \in [0,1]^K ig| i^*(oldsymbol{\lambda}) 
eq i ig\} \qquad ext{vs} \qquad \mathcal{H}_1 \ = \ \{oldsymbol{\mu}\}$$

We denote the relative entropy aka Kullback Leibler divergence from Bernoulli x to Bernoulli y by  $kl(x, y) = x \ln \frac{x}{y} + (1 - x) \ln \frac{1 - x}{1 - y}$ .

We denote the relative entropy aka Kullback Leibler divergence from Bernoulli x to Bernoulli y by  $kl(x, y) = x \ln \frac{x}{y} + (1 - x) \ln \frac{1 - x}{1 - y}$ .

Fix  $\mu$  and  $\lambda \in \neg i^*(\mu)$ 

#### Fact

Fix the number  $n_k$  of samples form each arm k. If the algorithm is  $\delta$ -correct at time  $\tau = \sum_{k=1}^{K} n_k$  for both  $\mu$ ,  $\lambda$ , then

$$\sum_{k=1}^{K} \mathsf{N}_k \, \mathsf{kl}(\mu_k,\lambda_k) \, \geq \, \, \mathsf{kl}(\delta,1-\delta)$$

We denote the relative entropy aka Kullback Leibler divergence from Bernoulli x to Bernoulli y by  $kl(x, y) = x \ln \frac{x}{y} + (1 - x) \ln \frac{1 - x}{1 - y}$ .

Fix  $\mu$  and  $\lambda \in \neg i^*(\mu)$ 

#### Fact

Fix the number  $n_k$  of samples form each arm k. If the algorithm is  $\delta$ -correct at time  $\tau = \sum_{k=1}^{K} n_k$  for both  $\mu$ ,  $\lambda$ , then

$$\sum_{k=1}^{K} \mathsf{N}_k \, \mathsf{kl}(\mu_k,\lambda_k) \, \geq \, \, \mathsf{kl}(\delta,1-\delta)$$

Proof by data compression, tensorisation of KL,  $\delta$ -PAC.

$$H := (I_1, X_1, \ldots, I_{\tau}, X_{\tau}, \hat{\imath})$$

We denote the relative entropy aka Kullback Leibler divergence from Bernoulli x to Bernoulli y by  $kl(x, y) = x \ln \frac{x}{y} + (1 - x) \ln \frac{1 - x}{1 - y}$ .

Fix  $\mu$  and  $\lambda \in \neg i^*(\mu)$ 

#### Fact

Fix the number  $n_k$  of samples form each arm k. If the algorithm is  $\delta$ -correct at time  $\tau = \sum_{k=1}^{K} n_k$  for both  $\mu$ ,  $\lambda$ , then

$$\sum_{k=1}^{K} \mathsf{N}_k \, \mathsf{kl}(\mu_k,\lambda_k) \geq \, \mathsf{kl}(\delta,1-\delta)$$

Proof by data compression, tensorisation of KL,  $\delta$ -PAC.

$$H := (I_1, X_1, \dots, I_{\tau}, X_{\tau}, \hat{\imath}) \qquad \mathsf{KL}\begin{pmatrix}\mathsf{alg}\\\mathbb{P}\\\mu\end{pmatrix} \begin{pmatrix}\mathsf{alg}\\\mathcal{X}\end{pmatrix} \geq \mathsf{kl}\begin{pmatrix}\mathsf{alg}\\\mathbb{P}\\\hat{\imath} \neq i^*(\mu) \end{pmatrix}, \overset{\mathsf{alg}}{\mathbb{P}}\{\hat{i} \neq i^*(\mu)\}, \end{pmatrix}$$

If we follow the sampling and stopping rules, we get on average  $\mathbb{E}_{\mu}[N_k(\tau)]$  samples from arm k at time  $\tau$ .

If we follow the sampling and stopping rules, we get on average  $\mathbb{E}_{\mu}[N_k(\tau)]$  samples from arm k at time  $\tau$ . With that,

$$\sum_{k=1}^{K} \mathbb{E}[N_k( au)] \, \mathsf{kl}(\mu_k, \lambda_k) \; \geq \; \mathsf{kl}(\delta, 1-\delta)$$

If we follow the sampling and stopping rules, we get on average  $\mathbb{E}_{\mu}[N_k(\tau)]$  samples from arm k at time  $\tau$ . With that,

$$\sum_{k=1}^{K} \mathop{\mathbb{E}}_{\mu}[\mathsf{N}_{k}( au)] \, \mathsf{kl}(\mu_{k},\lambda_{k}) \; \geq \; \mathsf{kl}(\delta,1-\delta)$$

Since this has to hold for all  $\lambda \in 
eg i^*(\mu)$ , we find

$$\inf_{\boldsymbol{\lambda}\in\neg i^{*}(\boldsymbol{\mu})}\sum_{k=1}^{K} \mathbb{E}[N_{k}(\tau)] \operatorname{kl}(\mu_{k},\lambda_{k}) \geq \operatorname{kl}(\delta,1-\delta)$$

If we follow the sampling and stopping rules, we get on average  $\mathbb{E}_{\mu}[N_k(\tau)]$  samples from arm k at time  $\tau$ . With that,

$$\sum_{k=1}^{K} \mathbb{E}[N_k( au)] \, \mathsf{kl}(\mu_k, \lambda_k) \, \geq \, \, \mathsf{kl}(\delta, 1-\delta)$$

Since this has to hold for all  $\lambda \in 
eg i^*(\mu)$ , we find

$$\inf_{\boldsymbol{\lambda}\in \neg i^*(\boldsymbol{\mu})} \sum_{k=1}^{\mathcal{K}} \mathop{\mathbb{E}}_{\boldsymbol{\mu}}[\textit{N}_k(\tau)] \, \mathsf{kl}(\mu_k,\lambda_k) \, \geq \, \, \mathsf{kl}(\delta,1-\delta)$$

Even for the best algorithm, we find

$$\max_{w \in \bigtriangleup_K} \inf_{\boldsymbol{\lambda} \in \neg i^*(\boldsymbol{\mu})} \mathbb{E}[\tau] \sum_{k=1}^K w_k \operatorname{kl}(\mu_k, \lambda_k) \geq \operatorname{kl}(\delta, 1-\delta)$$

# Lower Bound

Recall that to give answer i we need to reject the composite hypothesis

$$eg i := \{ \boldsymbol{\lambda} \in [0,1]^K \mid i^*(\boldsymbol{\lambda}) \neq i \}.$$

### Lower Bound

Recall that to give answer *i* we need to reject the composite hypothesis

 $eg i := \{ \boldsymbol{\lambda} \in [0,1]^{K} \mid i^{*}(\boldsymbol{\lambda}) \neq i \}.$ 

#### Theorem (Garivier and Kaufmann 2016)

Fix a  $\delta$ -correct strategy. Then for every bandit model  $\mu$ 

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

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

$$\frac{1}{T^*(\boldsymbol{\mu})} = \max_{\boldsymbol{w} \in \bigtriangleup_K} \min_{\boldsymbol{\lambda} \in \neg i^*(\boldsymbol{\mu})} \sum_{k=1}^K w_k \operatorname{kl}(\mu_k, \lambda_k).$$

# Example

Unknown environment







Arm









Algorithms

# Algorithm

Lower bound gives us more than just sample complexity.



## Algorithm



Lower bound gives us more than just sample complexity. We also get the oracle weight map

$$w^*(\mu) \ \coloneqq \ \max_{w \in riangle_K} \inf_{oldsymbol{\lambda} \in 
eg i^*(\mu)} \sum_{k=1}^K w_k \, \mathsf{kl}(\mu_k, \lambda_k)$$

What if **plug-in** empirical mean  $\hat{\mu}(t-1)$  and play arm  $I_t \sim w^* (\hat{\mu}(t-1))$ .

# Algorithm



Lower bound gives us more than just sample complexity. We also get the oracle weight map 47

$$w^*(\mu) \ \coloneqq \ \max_{w \in riangle_K} \inf_{oldsymbol{\lambda} \in 
eg i^*(\mu)} \sum_{k=1}^K w_k \, \mathsf{kl}(\mu_k, \lambda_k)$$

What if **plug-in** empirical mean  $\hat{\mu}(t-1)$  and play arm  $I_t \sim w^*(\hat{\mu}(t-1))$ .

#### Theorem (Degenne and Koolen, 2019)

Take set-valued interpretation of arg max defining  $w^*$ . Then  $\mu \mapsto w^*(\mu)$  is upper-hemicontinuous and convex-valued. Suitable tracking ensures that as  $\hat{\mu}(t) \to \mu$ , any sequence of choices  $w_t \in w^*(\hat{\mu}(t-1))$  has

$$\min_{oldsymbol{w}\inoldsymbol{w}^*(oldsymbol{\mu})} \left\|oldsymbol{w}_t - oldsymbol{w}
ight\|_{\infty} o \mathsf{0} \quad ext{as } t o \infty.$$

How to stop as soon as possible?

How to stop as soon as possible?

Need confidence region optimised for problem.

How to stop as soon as possible?

Need confidence region optimised for problem.

Say answer  $\hat{\imath}_t$  is correct for  $\hat{\mu}(t)$ , and we want to conclude that it is correct for  $\mu$ .

How to stop as soon as possible?

Need confidence region optimised for problem.

Say answer  $\hat{\imath}_t$  is correct for  $\hat{\mu}(t)$ , and we want to conclude that it is correct for  $\mu$ .

We need to reject the possibility that any other answer is correct, i.e. the composite hypothesis

$$\neg \hat{\imath}_t := \{ \boldsymbol{\lambda} \mid i^*(\lambda_k) \neq \hat{\imath}_t \}.$$

How to stop as soon as possible?

Need confidence region optimised for problem.

Say answer  $\hat{\imath}_t$  is correct for  $\hat{\mu}(t)$ , and we want to conclude that it is correct for  $\mu$ .

We need to reject the possibility that any other answer is correct, i.e. the composite hypothesis

$$eg \hat{\imath}_t \coloneqq \{ \boldsymbol{\lambda} \mid i^*(\lambda_k) \neq \hat{\imath}_t \}.$$

How to measure the evidence against  $\neg \hat{\imath}_t$ ?

Generalised Likelihood Ratio Test (GLRT) Statistic

$$\Lambda_t := \ln \frac{\max_{\lambda \notin \neg \hat{\imath}_t} p(data|\lambda)}{\max_{\lambda \in \neg \hat{\imath}_t} p(data|\lambda)} = \min_{\lambda \in \neg \hat{\imath}_t} \sum_{k=1}^K N_k(t) \operatorname{kl}(\hat{\mu}_{t,k}, \lambda_k)$$

One can show that

Theorem (Kaufmann and Koolen, 2021)

For every bandit  $\mu$ 

$$\mathbb{P}_{\mu}\left(\exists t: \underbrace{\Lambda_t \geq \ln \frac{1}{\delta} + O(\ln \ln t)}_{GLRT \text{ is big}} \text{ and } \underbrace{\hat{\imath}_t \neq i^*(\mu)}_{mistake}\right) \leq \delta.$$

Cool, so we can stop when  $\Lambda_t \geq \ln \frac{1}{\delta} + O(\ln \ln t)$  and safely output  $\hat{\imath}_t$
# Obtaining a sampling rule from GLRT

So we can stop when we see a big value in the GLRT

$$\Lambda_t = \min_{\boldsymbol{\lambda} \in \neg \hat{\imath}_t} \sum_{k=1}^{K} N_k(t) \operatorname{kl} \left( \hat{\mu}_{t,k}, \lambda_k \right)$$

# Obtaining a sampling rule from GLRT

So we can stop when we see a big value in the GLRT

$$\Lambda_t = \min_{oldsymbol{\lambda} \in 
eg \hat{\imath}_t} \sum_{k=1}^K N_k(t) \, \mathrm{kl} \left( \hat{\mu}_{t,k}, \lambda_k 
ight)$$

To stop early, the sampling rule should drive  $\Lambda_t$  up.

# Obtaining a sampling rule from GLRT

So we can stop when we see a big value in the GLRT

$$\Lambda_t = \min_{oldsymbol{\lambda} \in 
eg \hat{\imath}_t} \sum_{k=1}^K N_k(t) \operatorname{kl}\left(\hat{\mu}_{t,k}, \lambda_k\right)$$

To stop early, the sampling rule should drive  $\Lambda_t$  up.

The GLRT/round is maximised by sampling with proportions

$$w_t \coloneqq rgmax_{w \in riangle \kappa} \min_{oldsymbol{\lambda} \in 
eg \hat{\iota}_t} \; \sum_{k=1}^K w_k \, \mathsf{kl} \left( \hat{\mu}_{t,k}, \lambda_k 
ight)$$

Again empirical plug-in of oracle weights.

# Track and Stop Framework

**Definition (Track-and-Stop)** 

• Stop and output 
$$\hat{\imath}_t$$
 when  $\Lambda_t = \min_{\lambda \in \neg \hat{\imath}_t} \sum_{k=1}^K N_k(t) \operatorname{kl}(\hat{\mu}_{t,k}, \lambda_k) \ge \operatorname{ln} \frac{1}{\delta} + O(\operatorname{ln} \operatorname{ln} t)$   
• Sample with proportions  $w_t = \underset{w \in \Delta_K}{\operatorname{arg max}} \min_{\lambda \in \neg \hat{\imath}_t} \sum_{k=1}^K w_k \operatorname{kl}(\hat{\mu}_{t,k}, \lambda_k)$ 

## **Track and Stop Framework**

**Definition (Track-and-Stop)** 

• Stop and output 
$$\hat{\imath}_t$$
 when  $\Lambda_t = \min_{\lambda \in \neg \hat{\imath}_t} \sum_{k=1}^K N_k(t) \operatorname{kl}(\hat{\mu}_{t,k}, \lambda_k) \ge \operatorname{ln} \frac{1}{\delta} + O(\operatorname{ln} \operatorname{ln} t)$   
• Sample with proportions  $w_t = \underset{w \in \Delta_K}{\operatorname{arg\,max}} \min_{\lambda \in \neg \hat{\imath}_t} \sum_{k=1}^K w_k \operatorname{kl}(\hat{\mu}_{t,k}, \lambda_k)$ 

#### Theorem (Garivier and Kaufmann, 2016)

The expected stopping time on any bandit  $\mu$  is

$$\mathbb{E}[\tau] \leq \frac{\ln \frac{1}{\delta}}{\max_{\boldsymbol{w} \in \triangle_{K}} \min_{\boldsymbol{\lambda} \in \neg i^{*}(\boldsymbol{\mu})} \sum_{k=1}^{K} w_{k} \operatorname{kl}(\mu_{k}, \lambda_{k})} + small$$

and this matches lower bounds.

# Upshot of Track-and-Stop

The cost for learning the correct answer is

$$\frac{\ln \frac{1}{\delta}}{\max_{\boldsymbol{w} \in \bigtriangleup_{K}} \min_{\boldsymbol{\lambda} \in \neg i^{*}(\boldsymbol{\mu})} \sum_{k=1}^{K} w_{k} \operatorname{kl}(\mu_{k}, \lambda_{k})} + \operatorname{small}$$

Good:

- reliable
- practically efficient
- optimal (stop search for improvements)

#### Bad:

- Requires solving  $\arg \max_w \min_{\lambda} \cdots$
- I hid some details

# Discussion

Racing samples uniformly, and eliminates arms sequentially until one is left (Even-Dar, Mannor, and Mansour, 2006).

Can be generalised to MCTS (Teraoka, Hatano, and Takimoto, 2014).

Penultimate arm always sampled equally often as winner.

This by itself can result in a factor 2 sub-optimality in sample complexity.

In addition, how are we comparing arms?

Let's compare arm A with arm B, using samples, without knowing their qualities  $\mu_A$  or  $\mu_B$ .

Let's compare arm A with arm B, using samples, without knowing their qualities  $\mu_A$  or  $\mu_B$ .

#### Fact

If  $|\mu_A - \hat{\mu}_A| \le \epsilon_A$  and  $|\mu_B - \hat{\mu}_B| \le \epsilon_B$  then

$$ig|(\mu_{\mathcal{A}}-\mu_{\mathcal{B}})-(\hat{\mu}_{\mathcal{A}}-\hat{\mu}_{\mathcal{B}})ig| \ \le \ \epsilon_{\mathcal{A}}+\epsilon_{\mathcal{B}}$$

Let's compare arm A with arm B, using samples, without knowing their qualities  $\mu_A$  or  $\mu_B$ .

#### Fact

If 
$$|\mu_A - \hat{\mu}_A| \le \epsilon_A$$
 and  $|\mu_B - \hat{\mu}_B| \le \epsilon_B$  then

$$\left| \left( \mu_{\mathcal{A}} - \mu_{\mathcal{B}} 
ight) - \left( \hat{\mu}_{\mathcal{A}} - \hat{\mu}_{\mathcal{B}} 
ight) 
ight| ~\leq~ \epsilon_{\mathcal{A}} + \epsilon_{\mathcal{B}}$$



Let's compare arm A with arm B, using samples, without knowing their qualities  $\mu_A$  or  $\mu_B$ .

#### Fact

If 
$$|\mu_A - \hat{\mu}_A| \le \epsilon_A$$
 and  $|\mu_B - \hat{\mu}_B| \le \epsilon_B$  then

$$ig|(\mu_{\mathcal{A}}-\mu_{\mathcal{B}})-(\hat{\mu}_{\mathcal{A}}-\hat{\mu}_{\mathcal{B}})ig| ~\leq~ \epsilon_{\mathcal{A}}+\epsilon_{\mathcal{B}}$$



Let's compare arm A with arm B, using samples, without knowing their qualities  $\mu_A$  or  $\mu_B$ .

#### Fact

If 
$$|\mu_A - \hat{\mu}_A| \le \epsilon_A$$
 and  $|\mu_B - \hat{\mu}_B| \le \epsilon_B$  then

$$ig|(\mu_{A}-\mu_{B})-(\hat{\mu}_{A}-\hat{\mu}_{B})ig| ~\leq~ \epsilon_{A}+\epsilon_{B}$$



```
So far, we discussed active (\epsilon, \delta)-PAC learning for \epsilon = 0.
```

```
What about \epsilon > 0?
```

Lower bounds can be extended (Degenne and Koolen, 2019).

Yet algorithms more tricky: oracle weights  $w^*(\mu)$  now a **discontinuous** function of bandit instance  $\mu$ .

Simply ignoring that does not work. Terrible WSP.

Sticky-Track-and-Stop algorithm a "theoretical" fix.

Asymptotic optimality not quite so practical here. Also, only works for finitely many answers.

## Computation

In many cases the oracle weight map can be computed efficiently.

In some cases it cannot (Al Marjani and Proutiere, 2021; Ruitong, Ajallooeian, Szepesvári, and Müller, 2017).

In between are "theoretical" positive case. For example in Pareto Front Identification the run-time is  $O(\text{arms}^{\text{dimension}})$  (Crepon, Garivier, and Koolen, 2024).

Somehow every case is different.

Every classical CS problem becomes a pure exploration problem with bandit access to the inputs.

# Conclusion

# Conclusion

We saw

- How to decide which samples to collect
- Intuition
- Theory and algorithm design

We saw

- The specific question matters for the test!
- Lots to discuss and discover

# Let's talk!

# References i

- - Agrawal, S., W. M. Koolen, and S. Juneja (Dec. 2021). "Optimal Best-Arm Identification Methods for Tail-Risk Measures". In: Advances in Neural Information Processing Systems (NeurIPS) 34.
- Al Marjani, A. and A. Proutiere (2021). "Adaptive sampling for best policy identification in markov decision processes". In: International Conference on Machine Learning. PMLR.
- Chen, J., X. Chen, Q. Zhang, and Y. Zhou (2017). "Adaptive multiple-arm identification". In: International Conference on Machine Learning. PMLR.
- Chen, L., J. Li, and M. Qiao (2017). "Nearly instance optimal sample complexity bounds for top-k arm selection". In: *Artificial Intelligence and Statistics*. PMLR.
- Crepon, É., A. Garivier, and W. M. Koolen (Mar. 2024). "Sequential Learning of the Pareto Front for Multi-objective Bandits". To appear in AISTATS.
- Degenne, R. and W. M. Koolen (Dec. 2019). "Pure Exploration with Multiple Correct Answers". In: Advances in Neural Information Processing Systems (NeurIPS) 32.

# References ii

- Degenne, R., W. M. Koolen, and P. Ménard (Dec. 2019). "Non-Asymptotic Pure Exploration by Solving Games". In: Advances in Neural Information Processing Systems (NeurIPS) 32.
- Degenne, R., P. Ménard, X. Shang, and M. Valko (2020). "Gamification of pure exploration for linear bandits". In: International Conference on Machine Learning. PMLR.
- Even-Dar, E., S. Mannor, and Y. Mansour (2006). "Action Elimination and Stopping Conditions for the Multi-Armed Bandit and Reinforcement Learning Problems". In: Journal of Machine Learning Research 7.
- Garivier, A. and E. Kaufmann (2016). "Optimal Best arm Identification with Fixed Confidence". In: Proceedings of the 29th Conference On Learning Theory (COLT).
- Garivier, A., E. Kaufmann, and W. M. Koolen (June 2016). "Maximin Action Identification: A New Bandit Framework for Games". In: Proceedings of the 29th Annual Conference on Learning Theory (COLT).
- Garivier, A., P. Ménard, and L. Rossi (2017). "Thresholding bandit for dose-ranging: The impact of monotonicity". In: *arXiv preprint arXiv:1711.04454*.

# References iii

- Haddenhorst, B., V. Bengs, and E. Hüllermeier (2021). "Identification of the generalized Condorcet winner in multi-dueling bandits". In: Advances in Neural Information Processing Systems 34.
- Hao, B., T. Lattimore, and C. Szepesvari (2020). "Adaptive exploration in linear contextual bandit". In: International Conference on Artificial Intelligence and Statistics. PMLR.
- Jedra, Y. and A. Proutiere (2020). "Optimal best-arm identification in linear bandits". In: Advances in Neural Information Processing Systems 33.
- **i** Katariya, S., A. Tripathy, and R. Nowak (2019). **"Maxgap bandit: Adaptive algorithms for approximate ranking".** In: *Advances in Neural Information Processing Systems* 32.
- Kaufmann, E. and W. M. Koolen (Nov. 2021). "Mixture Martingales Revisited with Applications to Sequential Tests and Confidence Intervals". In: Journal of Machine Learning Research 22.246.
- Kaufmann, E., W. M. Koolen, and A. Garivier (Dec. 2018). "Sequential Test for the Lowest Mean: From Thompson to Murphy Sampling". In: Advances in Neural Information Processing Systems (NeurIPS) 31.

## References iv

- Kocák, T. and A. Garivier (2021). "Epsilon Best Arm Identification in Spectral Bandits.". In: IJCAI.
- Kone, C., E. Kaufmann, and L. Richert (2023). "Adaptive Algorithms for Relaxed **Pareto Set Identification**". In: *arXiv preprint arXiv:2307.00424*.
- Moulos, V. (2019). "Optimal best Markovian arm identification with fixed confidence". In: Advances in Neural Information Processing Systems 32.
- Ruitong, H., M. M. Ajallooeian, C. Szepesvári, and M. Müller (2017). "Structured Best Arm Identification with Fixed Confidence". In: Algorithmic Learning Theorem (ALT).
- Russac, Y., C. Katsimerou, D. Bohle, O. Cappé, A. Garivier, and W. M. Koolen (Dec. 2021). "A/B/n Testing with Control in the Presence of Subpopulations". In:

2021). A/B/II resting with control in the Presence of Subpopulations

Advances in Neural Information Processing Systems (NeurIPS) 34.

Teraoka, K., K. Hatano, and E. Takimoto (2014). "Efficient Sampling Method for Monte Carlo Tree Search Problem". In: IEICE Transactions 97-D.3.

#### References v

Tirinzoni, A., M. Pirotta, M. Restelli, and A. Lazaric (2020). **"An asymptotically optimal primal-dual incremental algorithm for contextual linear bandits".** In: *Advances in Neural Information Processing Systems* 33.