Instance Optimal Pure Exploration in Bandit Models



Wouter M. Koolen





Young European Queueing Theorists Workshop EURANDOM Wednesday 2nd November, 2022 Emilie Kaufmann

Rémy Degenne

Pierre Menard

Han Shao

Aurelién Garivier

Sandeep Juneja

Shubhada Agrawal

Christina Katsimerou

Olivier Cappé

Yoan Russac



Grand Goal: Interactive Machine Learning





Grand Goal: Interactive Machine Learning





Main scientific questions

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

Pure Exploration focuses on the statistical problem (learn the truth), while Reinforcement Learning focuses on behaviour (maximise reward).

Pure Exploration focuses on the statistical problem (learn the truth), while Reinforcement Learning focuses on behaviour (maximise reward).

Pure Exploration occurs as **sub-module** in some RL algorithms (i.e. Phased Q-Learning by Even-Dar, Mannor, and Mansour, 2002)

Pure Exploration focuses on the statistical problem (learn the truth), while Reinforcement Learning focuses on behaviour (maximise reward).

Pure Exploration occurs as **sub-module** in some RL algorithms (i.e. Phased Q-Learning by Even-Dar, Mannor, and Mansour, 2002)

Some problems approached with RL are in fact **better modelled** as pure exploration problems. Most notably MCTS for playing games.



1. A Taste of the Problem Space

- 2. Formal Setup
- 3. Sample Complexity Lower Bounds: Information Theory
- 4. Numerical Illustration of Characteristic Time and Oracle Proportions
- 5. Design of Algorithms: Equilibria
- 6. Conclusion

Stochastic Bandit

Experiments





Stochastic Bandit

Experiments



Outcomes

Instance (Unknown)

$$\mathbb{P}\left(\bigotimes \middle| \overleftrightarrow{\bigotimes} \right) = 1/6$$
$$\mathbb{P}\left(\bigotimes \middle| \overleftrightarrow{\odot} \right) = 4/6$$
$$\mathbb{P}\left(\bigotimes \middle| \overleftrightarrow{\odot} \right) = 3/6$$



$$\mathbb{P}\left(\bigotimes|\bigotimes]\right) = 1/6$$
$$\mathbb{P}\left(\bigotimes|\bigotimes]\right) = 4/6$$
$$\mathbb{P}\left(\bigotimes|\bigotimes]\right) = 3/6$$











Identification Problems

Problem (Even-Dar, Mannor, and Mansour, 2002)

Which arm has the highest mean

Arms: Bernoulli, Exp. Fam, bounded support, sub-Gaussian, ...

Problem (Yu and Nikolova, 2013)

Which arm has the highest α -quantile **Arms**: Unrestricted (on \mathbb{R})

Problem (Yu and Nikolova, 2013)

Which arm has the smallest Conditional Value at Risk. Arms: Exp. Fam (trivial), bounded $(1 + \epsilon)^{\text{th}}$ moment



Set-valued Queries

Top-M, Thresholding, All ϵ -optimal, all-better-than-control **Structure**: standard bandit

Ranking-related queries: Borda/Condorcet winner (Yue et al., 2012)

Structure: duelling bandit

Best arm in stratified population (Russac et al., 2021)

Structure: contextual bandit

Minimax action in extensive form game tree (Teraoka, Hatano, and Takimoto, 2014) Structure: game tree with stochastic leaves

Shortest Path (Chen et al., 2014)

Structure: graph with stochastic edge costs



1. A Taste of the Problem Space

2. Formal Setup

- 3. Sample Complexity Lower Bounds: Information Theory
- 4. Numerical Illustration of Characteristic Time and Oracle Proportions
- 5. Design of Algorithms: Equilibria
- 6. Conclusion

Assumption: Bernoulli Multi-Armed Bandit

K Bernoulli arms with unknown means $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K) \in [0, 1]^K$.

Assumption: Bernoulli Multi-Armed Bandit

K Bernoulli arms with unknown means $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K) \in [0, 1]^K$.

The **best** arm is

 $i^*(\mu) = \arg \max_i \mu_i$

Assumption: Bernoulli Multi-Armed Bandit

K Bernoulli arms with unknown means $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K) \in [0, 1]^K$.

The **best** arm is

 $i^*(\mu) = \arg \max_i \mu_i$

BAI-MAB Protocol

for $t = 1, 2, \dots$ until Learner decides to stop

- Learner picks arm $A_t \in [K]$
- Learner observes $X_t \sim \text{Bernoulli}(\mu_{A_t})$

Learner recommends $\hat{l} \in [K]$.

BAI-MAB Protocol

for $t = 1, 2, \dots$ until Learner decides to stop

- Learner picks arm $A_t \in [K]$
- Learner observes $X_t \sim \text{Bernoulli}(\mu_{A_t})$

Learner recommends $\hat{l} \in [K]$.

Let $\tau \in \mathbb{N} \cup \{\infty\}$ denote the # rounds after which Learner stops.

BAI-MAB Protocol

for $t = 1, 2, \dots$ until Learner decides to stop

- Learner picks arm $A_t \in [K]$
- Learner observes $X_t \sim \text{Bernoulli}(\mu_{A_t})$

Learner recommends $\hat{l} \in [K]$.

Let $\tau \in \mathbb{N} \cup \{\infty\}$ denote the # rounds after which Learner stops.

Definition

Learner is
$$\delta$$
-PAC if $\mathbb{P}_{\mu}\left\{\underbrace{\tau < \infty \text{ and } \hat{l} \neq i^{*}(\mu)}_{\text{a mistake}}\right\} \leq \delta$ for all $\mu \in [0, 1]^{K}$.

BAI-MAB Protocol

for $t = 1, 2, \dots$ until Learner decides to stop

- Learner picks arm $A_t \in [K]$
- Learner observes $X_t \sim \text{Bernoulli}(\mu_{A_t})$

Learner recommends $\hat{l} \in [K]$.

Let $\tau \in \mathbb{N} \cup \{\infty\}$ denote the # rounds after which Learner stops.

Definition

Learner is
$$\delta$$
-**PAC** if $\mathbb{P}_{\mu}\left\{\underbrace{\tau < \infty \text{ and } \hat{l} \neq i^{*}(\mu)}_{\text{a mistake}}\right\} \leq \delta$ for all $\mu \in [0, 1]^{K}$.

Definition

We call $\mathbb{E}_{\mu}[\tau]$ the sample complexity of Learner in bandit μ .

BAI-MAB Protocol

for $t = 1, 2, \dots$ until Learner decides to stop

- Learner picks arm $A_t \in [K]$
- Learner observes $X_t \sim \text{Bernoulli}(\mu_{A_t})$

Learner recommends $\hat{l} \in [K]$.

Let $\tau \in \mathbb{N} \cup \{\infty\}$ denote the # rounds after which Learner stops.

Definition

Learner is
$$\delta$$
-**PAC** if $\mathbb{P}_{\mu}\left\{\underbrace{\tau < \infty \text{ and } \hat{l} \neq i^{*}(\mu)}_{\text{a mistake}}\right\} \leq \delta$ for all $\mu \in [0, 1]^{K}$.

Definition

We call $\mathbb{E}_{\mu}[\tau]$ the sample complexity of Learner in bandit μ .

Goal: computationally efficient δ -PAC algorithms with minimal sample complexity.



- 1. A Taste of the Problem Space
- 2. Formal Setup

3. Sample Complexity Lower Bounds: Information Theory

- 4. Numerical Illustration of Characteristic Time and Oracle Proportions
- 5. Design of Algorithms: Equilibria
- 6. Conclusion

Intuition, going back at least to Lai and Robbins (1985)

(Spectacular) difference in behaviour must be due to (spectacular) difference in observations.

Being δ -PAC on μ and λ with $i^*(\mu) \neq i^*(\lambda)$ requires gathering enough discriminating data.

Intuition, going back at least to Lai and Robbins (1985)

(Spectacular) difference in behaviour must be due to (spectacular) difference in observations.

Being δ -PAC on μ and λ with $i^*(\mu) \neq i^*(\lambda)$ requires gathering enough discriminating data.

$$\begin{array}{c} \mathbb{P}\left(\textcircled{\textcircled{\sc oneq}}\right) &=& 1/6\\ \mathbb{P}\left(\textcircled{\textcircled{\sc oneq}}\right) &=& 4/6\\ \mathbb{P}\left(\textcircled{\textcircled{\sc oneq}}\right) &=& 3/6 \end{array} \leftarrow i^* \qquad \begin{array}{c} \mathbb{P}\left(\textcircled{\textcircled{\sc oneq}}\right) &=& 1/4\\ \mathbb{P}\left(\textcircled{\textcircled{\sc oneq}}\right) &=& 2/4\\ \mathbb{P}\left(\textcircled{\textcircled{\sc oneq}}\right) &=& 3/4 \end{array} \leftarrow i^* \end{array}$$

If δ -PAC algorithm samples t rounds with arm frequencies $(\frac{1}{6}, \frac{3}{6}, \frac{2}{6})$, then

$$t\,\frac{1}{6}\,\mathsf{KL}\left(\frac{1}{6},\frac{1}{4}\right)+t\,\frac{3}{6}\,\mathsf{KL}\left(\frac{4}{6},\frac{2}{4}\right)+t\,\frac{2}{6}\,\mathsf{KL}\left(\frac{3}{6},\frac{3}{4}\right) \ \geq \ \mathsf{KL}(\delta,1-\delta)\approx\ln\frac{1}{\delta}$$

Intuition, going back at least to Lai and Robbins (1985)

(Spectacular) difference in behaviour must be due to (spectacular) difference in observations.

Being δ -PAC on μ and λ with $i^*(\mu) \neq i^*(\lambda)$ requires gathering enough discriminating data.

$$\begin{array}{c} \mathbb{P}\left(\textcircled{\textcircled{\sc oneq}}\right) &=& 1/6\\ \mathbb{P}\left(\textcircled{\textcircled{\sc oneq}}\right) &=& 4/6\\ \mathbb{P}\left(\textcircled{\textcircled{\sc oneq}}\right) &=& 3/6 \end{array} \leftarrow i^* \qquad \begin{array}{c} \mathbb{P}\left(\textcircled{\textcircled{\sc oneq}}\right) &=& 1/4\\ \mathbb{P}\left(\textcircled{\textcircled{\sc oneq}}\right) &=& 2/4\\ \mathbb{P}\left(\textcircled{\textcircled{\sc oneq}}\right) &=& 3/4 \end{array} \leftarrow i^* \end{array}$$

If δ -PAC algorithm samples t rounds with arm frequencies $(\frac{1}{6}, \frac{3}{6}, \frac{2}{6})$, then

$$t\,\frac{1}{6}\,\mathsf{KL}\left(\frac{1}{6},\frac{1}{4}\right)+t\,\frac{3}{6}\,\mathsf{KL}\left(\frac{4}{6},\frac{2}{4}\right)+t\,\frac{2}{6}\,\mathsf{KL}\left(\frac{3}{6},\frac{3}{4}\right) \ \geq \ \mathsf{KL}(\delta,1-\delta)\approx\ln\frac{1}{\delta}$$

At typical $\delta = 0.1$: 0.0796 $t \ge 1.757$ $t \ge \frac{1.757}{0.0796} = 28.9$

Intuition, going back at least to Lai and Robbins (1985)

(Spectacular) difference in behaviour must be due to (spectacular) difference in observations.

Being δ -PAC on μ and λ with $i^*(\mu) \neq i^*(\lambda)$ requires gathering enough discriminating data.

Intuition, going back at least to Lai and Robbins (1985)

(Spectacular) difference in behaviour must be due to (spectacular) difference in observations.

Being δ -PAC on μ and λ with $i^*(\mu) \neq i^*(\lambda)$ requires gathering enough discriminating data.

Define the **alternative** to μ by $Alt(\mu) \coloneqq \{bandit \ \lambda | i^*(\lambda) \neq i^*(\mu)\}.$

Intuition, going back at least to Lai and Robbins (1985)

(Spectacular) difference in behaviour must be due to (spectacular) difference in observations.

Being δ -PAC on μ and λ with $i^*(\mu) \neq i^*(\lambda)$ requires gathering enough discriminating data.

Define the alternative to μ by Alt $(\mu) \coloneqq \{$ bandit $\lambda | i^*(\lambda) \neq i^*(\mu) \}$.

Theorem (Castro 2014; Garivier and Kaufmann 2016) *Fix a* δ *-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

$$\frac{1}{T^*(\mu)} = \underbrace{\max_{w \in \triangle_K} \min_{\lambda \in Alt(\mu)} \sum_{i=1}^K w_i \operatorname{KL}(\mu_i, \lambda_i)}_{w_i \in Alt(\mu)}$$

Sample complexity lower bound at bandit μ governed by characteristic time

$$\mathbb{E}_{\boldsymbol{\mu}}[\tau] \geq \left[\max_{\boldsymbol{w} \in \bigtriangleup_{\kappa}} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^{\kappa} w_i \operatorname{\mathsf{KL}}(\mu_i, \lambda_i) \right]^{-1} \ln \frac{1}{\delta}$$

Sample complexity lower bound at bandit μ governed by characteristic time

$$\mathbb{E}_{\boldsymbol{\mu}}[\tau] \geq \left[\max_{\boldsymbol{w} \in \bigtriangleup_{K}} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^{K} w_{i} \operatorname{\mathsf{KL}}(\mu_{i}, \lambda_{i}) \right]^{-1} \ln \frac{1}{\delta}$$

Moreover, matching algorithms must sample arms with oracle proportions

$$oldsymbol{w}^*(oldsymbol{\mu}) \ \coloneqq \ rgge w \in riangle_{\kappa} \min_{oldsymbol{\lambda} \in \mathsf{Alt}(oldsymbol{\mu})} \ \sum_{i=1}^K w_i \, \mathsf{KL}(\mu_i, \lambda_i)$$



- 1. A Taste of the Problem Space
- 2. Formal Setup
- 3. Sample Complexity Lower Bounds: Information Theory

4. Numerical Illustration of Characteristic Time and Oracle Proportions

- 5. Design of Algorithms: Equilibria
- 6. Conclusion

Examples: variations of Best Arm question



- Sample complexities vastly different between questions
- Optimal allocation depends strongly on the specific question being asked

Best Arm Identification (BAI)

$$i^*(\mu) \coloneqq rg\max_{a \in \mathcal{A}} \mu_a$$



where
$$\mathcal{A} = \{A, B, C, D\}$$



All-Better-than-the-Control (ABC)

$$i^*(\boldsymbol{\mu}) \coloneqq \left\{ \boldsymbol{a} \in \{B, C, D\} \mid \mu_{\boldsymbol{a}} \geq \mu_{\boldsymbol{A}} \right\}$$



$$\left(i^{*}(\boldsymbol{\mu}) \ \coloneqq \ \left\{ \boldsymbol{a} \in \mathcal{A} \mid \mu_{\boldsymbol{a}} \geq \gamma
ight\}
ight)$$



Top-2



$$i^*(\mu) \coloneqq \left\{ a \in \mathcal{A} \mid \mu_a \geq \mu^* - \epsilon
ight\}$$
 where $\mu^* = \max_{a \in \mathcal{A}} \mu_a$



 $i^*(oldsymbol{\mu}) \ \coloneqq \ ext{arg max}ig\{ ext{max}ig\{\mu_A,\mu_Big\}, ext{max}ig\{\mu_C,\mu_Dig\}ig\}$



 $i^*(\boldsymbol{\mu}) \coloneqq \operatorname{arg\,max} \{ \min \{ \mu_A, \mu_B \}, \min \{ \mu_C, \mu_D \} \}$



$$i^*(\boldsymbol{\mu}) \coloneqq \operatorname{arg\,max} \{\mu_{\boldsymbol{A}} - \mu_{\boldsymbol{B}}, \mu_{\boldsymbol{C}} - \mu_{\boldsymbol{D}}\}$$



Tree Search Example: Backward Induction Computation



Tree Search Example: Backward Induction Computation



Tree Search Example: Best action at root

$$i^*(\mu) := \arg \max_i \min_j \max_k \min_l \mu_{ijkl} \in \{\text{left, right}\}$$



Overview of Optimal Sampling Allocations



- Sample complexities vastly different between questions
- Optimal allocation depends strongly on the specific question being asked



- 1. A Taste of the Problem Space
- 2. Formal Setup
- 3. Sample Complexity Lower Bounds: Information Theory
- 4. Numerical Illustration of Characteristic Time and Oracle Proportions
- 5. Design of Algorithms: Equilibria
- 6. Conclusion

$$\max_{\boldsymbol{w} \in \bigtriangleup_{\kappa}} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^{\kappa} w_i \operatorname{\mathsf{KL}}(\mu_i, \lambda_i)$$

$$\max_{\boldsymbol{w} \in \bigtriangleup_{\kappa}} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^{\kappa} w_i \operatorname{\mathsf{KL}}(\mu_i, \lambda_i)$$

Matching algorithms must sample arms with argmax proportions $w^*(\mu)$.

$$\max_{\boldsymbol{w} \in \bigtriangleup_{\kappa}} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^{\kappa} w_i \operatorname{\mathsf{KL}}(\mu_i, \lambda_i)$$

Matching algorithms must sample arms with argmax proportions $w^*(\mu)$.

Main issue: Bandit instance μ unknown

$$\max_{\boldsymbol{w} \in \bigtriangleup_{\kappa}} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^{\kappa} w_i \operatorname{\mathsf{KL}}(\mu_i, \lambda_i)$$

Matching algorithms must sample arms with argmax proportions $w^*(\mu)$.

Main issue: Bandit instance μ unknown

Approach: plug in estimate $\hat{\mu}_t$ (Garivier and Kaufmann, 2016)

$$\max_{\boldsymbol{w} \in \bigtriangleup_{K}} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^{K} w_{i} \operatorname{\mathsf{KL}}(\mu_{i}, \lambda_{i})$$

online learning

Approx. solve saddle point problem iteratively: $w_1, w_2, \ldots o w^*(\mu)$

$$\max_{\boldsymbol{w} \in \triangle_K} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^K w_i \, \mathsf{KL}(\mu_i, \lambda_i)$$

Import

learning

Approx. solve saddle point problem iteratively: $w_1, w_2, \ldots \to w^*(\mu)$ Main pipeline (Degenne, Koolen, and Ménard, 2019):

- Pick arm $A_t \sim w_t$
- Plug-in estimate $\hat{\mu}_t$ (so problem is **shifting**).
- Advance the saddle point solver one iteration per bandit interaction.
- Add optimism to gradients to induce exploration $(\hat{\mu}_t
 ightarrow \mu)$.
- Compose regret bound, concentration and optimism to get finite-confidence guarantee.

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

Import

Chniques ning

Approx. solve saddle point problem iteratively: $w_1, w_2, \ldots \rightarrow w^*(\mu)$ Main pipeline (Degenne, Koolen, and Ménard, 2019):

- Pick arm $A_t \sim w_t$
- Plug-in estimate $\hat{\mu}_t$ (so problem is **shifting**).
- Advance the saddle point solver one iteration per bandit interaction.
- Add optimism to gradients to induce exploration $(\hat{\mu}_t
 ightarrow \mu).$
- Compose regret bound, concentration and optimism to get finite-confidence guarantee.

Theorem (Instance-Optimality)

For every $\delta \in (0,1)$, the sample complexity is bounded by $\mathbb{E}_{\mu}[\tau] \leq T^*(\mu) \ln \frac{1}{\delta} + o(\ln \frac{1}{\delta})$.

Sketch of Argument

As long as we do not stop (and concentration holds):

$$\ln \frac{1}{\delta} \geq \inf_{\lambda \in \operatorname{Alt}(\hat{\mu}_{t})} \sum_{k=1}^{K} N_{t}^{k} \operatorname{KL}(\mu^{k}, \lambda^{k}) \qquad (\text{stop rule})$$

$$\approx \inf_{\lambda \in \operatorname{Alt}(\mu)} \sum_{s=1}^{t} \sum_{k=1}^{K} w_{s}^{k} \operatorname{KL}(\mu^{k}, \lambda^{k}) \qquad (\text{tracking})$$

$$\geq \sum_{s=1}^{t} \sum_{k=1}^{K} w_{s}^{k} \mathbb{E}_{\lambda \sim q_{s}} \operatorname{KL}(\mu^{k}, \lambda^{k}) - R_{t}^{\lambda} \qquad (\text{regret } \lambda)$$

$$\geq \max_{k} \sum_{s=1}^{t} \mathbb{E}_{\lambda \sim q_{s}} \operatorname{KL}(\mu^{k}, \lambda^{k}) - R_{t}^{\lambda} - R_{t}^{k} \qquad (\text{regret } k)$$

$$\geq t \inf_{q \in \mathcal{P}(\operatorname{Alt}(\mu))} \max_{k} \mathbb{E}_{\lambda \sim q} \operatorname{KL}(\mu^{k}, \lambda^{k}) - O(\sqrt{t})$$

$$= t \underbrace{\max_{w \in \Delta_{K}} \min_{\lambda \in \operatorname{Alt}(\mu)}} \sum_{i=1}^{K} w_{i} \operatorname{KL}(\mu_{i}, \lambda_{i}) - O(\sqrt{t})$$



- 1. A Taste of the Problem Space
- 2. Formal Setup
- 3. Sample Complexity Lower Bounds: Information Theory
- 4. Numerical Illustration of Characteristic Time and Oracle Proportions
- 5. Design of Algorithms: Equilibria

Canonical Path to Instance Optimality

- State-of-the-art performance in practise (some problems)
 - Best Arm Identification
 - All-better-than-Control
 - Minimax Game Tree Search

Canonical Path to Instance Optimality

- State-of-the-art performance in practise (some problems)
 - Best Arm Identification
 - All-better-than-Control
 - Minimax Game Tree Search
- Different ("fresh") structure compared to other techniques (confidence intervals, elimination, Thompson sampling, ...)

Canonical Path to Instance Optimality

- State-of-the-art performance in practise (some problems)
 - Best Arm Identification
 - All-better-than-Control
 - Minimax Game Tree Search
- Different ("fresh") structure compared to other techniques (confidence intervals, elimination, Thompson sampling, ...)
- Reduces identification problems to online learning (efficiently computing gradients/best response).

Canonical Path to Instance Optimality

- State-of-the-art performance in practise (some problems)
 - Best Arm Identification
 - All-better-than-Control
 - Minimax Game Tree Search
- Different ("fresh") structure compared to other techniques (confidence intervals, elimination, Thompson sampling, ...)
- Reduces identification problems to online learning (efficiently computing gradients/best response).
- Foundation for
 - Linear bandits
 - Contextual bandits
 - Optimal policy learning (reinforcement learning)

Wish list:

• Instance optimality for ($\epsilon,\delta)\text{-case}$ currently deeply asymptotic

Thanks!

References i

- Castro, R. M. (Nov. 2014). "Adaptive sensing performance lower bounds for sparse signal detection and support estimation". In: Bernoulli 20.4, pp. 2217–2246.
- Chen, S., T. Lin, I. King, M. Lyu, and W. Chen (2014). "Combinatorial Pure Exploration of Multi-Armed Bandits". In: Advances in Neural Information Processing Systems.
- 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. Ed. by H. Wallach, H. Larochelle, A. Beygelzimer,

F. d'Alché-Buc, E. Fox, and R. Garnett. Curran Associates, Inc., pp. 14492–14501.

 Even-Dar, E., S. Mannor, and Y. Mansour (2002). "PAC Bounds for Multi-armed Bandit and Markov Decision Processes". In: Computational Learning Theory, 15th Annual Conference on Computational Learning Theory, COLT 2002, Sydney, Australia, July 8-10, 2002, Proceedings. Ed. by J. Kivinen and R. H. Sloan. Vol. 2375. Lecture Notes in Computer Science. Springer, pp. 255–270.

References ii

- Garivier, A. and E. Kaufmann (2016). "Optimal Best arm Identification with Fixed Confidence". In: Proceedings of the 29th Conference On Learning Theory (COLT).
- Lai, T. L. and H. Robbins (1985). "Asymptotically efficient adaptive allocation rules". In:
 Advances in Applied Mathematics 6.1, pp. 4–22.
- 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: 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, pp. 392–398.
- Yu, J. Y. and E. Nikolova (2013). "Sample complexity of risk-averse bandit-arm selection". In: Twenty-Third International Joint Conference on Artificial Intelligence.
- Yue, Y., J. Broder, R. Kleinberg, and T. Joachims (2012). "The K-armed dueling bandits problem". In: Journal of Computer and System Sciences, Special Issue for COLT 2009, pp. 1538–1566.