Reinforcement Learning 2021 Guest Lecture on Pure Exploration



Wouter M. Koolen

Download these slides from https://wouterkoolen.info/Talks/2021-10-19.pdf!

Pure Exploration:

- PAC Learning
- Best Arm Identification
- Minimax Strategies in Noisy Games (Zero-Sum, Extensive Form)

Introduction and Motivation

Grand Goal: Interactive Machine Learning



Grand Goal: Interactive Machine Learning



Main scientific questions

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

Both are about *learning in uncertain environments*.

Both are about *learning in uncertain environments*.

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

Both are about *learning in uncertain environments*.

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)

Both are about *learning in uncertain environments*.

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.

Best Arm Identification

max

Environment (Multi-armed bandit model)

K distributions parameterised by their means $\mu = (\mu_1, \dots, \mu_K)$. The *best arm* is

$$i^* = \underset{i \in [\mathcal{K}]}{\operatorname{argmax}} \mu_i$$

max

Environment (Multi-armed bandit model)

K distributions parameterised by their means $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K).$ The *best arm* is

$$i^* = \underset{i \in [K]}{\operatorname{argmax}} \mu_i$$

Strategy

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

max

Environment (Multi-armed bandit model)

K distributions parameterised by their means $\mu = (\mu_1, \dots, \mu_K)$. The *best arm* is

$$i^* = rgmax \mu_i$$

 $i \in [K]$

Strategy

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

Realisation of interaction: $(I_1, X_1), \ldots, (I_{\tau}, X_{\tau}), \hat{I}$.

max

Environment (Multi-armed bandit model)

K distributions parameterised by their means $\mu = (\mu_1, \dots, \mu_K)$. The *best arm* is

$$i^* = rgmax \mu_i$$

 $i \in [K]$

Strategy

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

Realisation of interaction: $(I_1, X_1), \ldots, (I_{\tau}, X_{\tau}), \hat{I}$.

Two objectives: sample efficiency τ and correctness $\hat{l} = i^*$.

Objective

0

On bandit μ , strategy $(au, (I_t)_t, \hat{I})$ has

- error probability $\mathbb{P}_{m{\mu}}ig(\hat{l}
 eq i^*(m{\mu})ig)$, and
- sample complexity $\mathbb{E}_{\mu}[\tau]$.

Idea: constrain one, optimise the other.

Objective



On bandit μ , strategy $(au, (I_t)_t, \hat{I})$ has

- error probability $\mathbb{P}_{m{\mu}}ig(\hat{l}
 eq i^*(m{\mu})ig)$, and
- sample complexity $\mathbb{E}_{\mu}[\tau]$.

Idea: constrain one, optimise the other.

Definition

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

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

(Generalisation: output ϵ -best arm)

Objective



On bandit μ , strategy $(au, (I_t)_t, \hat{I})$ has

- error probability $\mathbb{P}_{m{\mu}}ig(\hat{l}
 eq i^*(m{\mu})ig)$, and
- sample complexity $\mathbb{E}_{\mu}[\tau]$.

Idea: constrain one, optimise the other.

Definition

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

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

(Generalisation: output ϵ -best arm)

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

Algorithms



► Sampling rule *I*_t?

Stopping rule τ ?

• Recommendation rule \hat{I} ?

$$\hat{I} = \operatorname{argmax}_{i \in [K]} \hat{\mu}_i(au)$$

where $\hat{\mu}(t)$ is empirical mean.

Algorithms



► Sampling rule *I*_t?

Stopping rule τ ?

• Recommendation rule \hat{I} ?

$$\hat{l} = \underset{i \in [K]}{\operatorname{argmax}} \hat{\mu}_i(\tau)$$

where $\hat{\mu}(t)$ is empirical mean.

Approach: start investigating lower bounds

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

Define the alternatives to μ by $\mathsf{Alt}(\mu) = \{\lambda | i^*(\lambda) \neq i^*(\mu)\}.$

Theorem (Castro 2014; Garivier and Kaufmann 2016)

Fix a δ -correct strategy. Then for every bandit model μ

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

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

$$\frac{1}{\mathcal{T}^*(\boldsymbol{\mu})} = \max_{\boldsymbol{w} \in \bigtriangleup_K} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^K w_i \, \mathsf{KL}(\mu_i \| \lambda_i).$$

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

Theorem (Castro 2014; Garivier and Kaufmann 2016)

Fix a δ -correct strategy. Then for every bandit model μ

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

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

$$\frac{1}{\mathcal{T}^*(\boldsymbol{\mu})} = \max_{\boldsymbol{w} \in \Delta_K} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^K w_i \, \mathsf{KL}(\mu_i \| \lambda_i).$$

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

Blackboard proof

Example



K = 5 arms, Bernoulli $\mu = (0.0, 0.1, 0.2, 0.3, 0.4)$.

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

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

Sampling Rule

Look at the lower bound again. Any good algorithm **must** sample with optimal (*oracle*) proportions

$$m{w}^*(m{\mu}) \ = \ rgmax_{m{w}\in riangle _{m{\kappa}}} \min_{m{\lambda}\in riangle _{m{k}}} \ \sum_{i=1}^K w_i \, { extsf{KL}}(\mu_i \| \lambda_i)$$

Sampling Rule

Look at the lower bound again. Any good algorithm **must** sample with optimal (*oracle*) proportions

$$m{w}^*(m{\mu}) \;=\; rgmax \min_{m{w}\in riangle _K} \; \min_{m{\lambda}\in riangle riangle (m{\mu})} \; \sum_{i=1}^K w_i \, \mathsf{KL}(\mu_i \| \lambda_i)$$

Track-and-Stop

Idea: draw $I_t \sim w^*(\hat{\mu}(t-1)).$

- ▶ Ensure $\hat{\mu}(t) \rightarrow \mu$ hence $N_i(t)/t \rightarrow w_i^*$ by "forced exploration"
- Draw arm with $N_i(t)/t$ below w_i^* (tracking)
- Computation of w* (reduction to 1d line search)

When can we stop?

When can we stop? When can we stop and give answer $\hat{\imath}$?

When can we stop? When can we stop and give answer $\hat{\imath}$? There is no plausible bandit model λ on which $\hat{\imath}$ is wrong.

When can we stop? When can we stop and give answer $\hat{\imath}$? There is no plausible bandit model λ on which $\hat{\imath}$ is wrong.

Definition

Generalized Likelihood Ratio (GLR) measure of evidence

$$\mathsf{GLR}_n(\hat{\imath}) := \ln \frac{\sup_{\mu:\hat{\imath} \in i^*(\mu)} P\left(X^n | A^n, \mu\right)}{\sup_{\lambda:\hat{\imath} \notin i^*(\lambda)} P\left(X^n | A^n, \lambda\right)}$$

When can we stop? When can we stop and give answer $\hat{\imath}$? There is no plausible bandit model λ on which $\hat{\imath}$ is wrong.

Definition

Generalized Likelihood Ratio (GLR) measure of evidence

$$\mathsf{GLR}_n(\hat{\imath}) := \ln \frac{\sup_{\mu:\hat{\imath} \in i^*(\mu)} P(X^n | A^n, \mu)}{\sup_{\lambda:\hat{\imath} \notin i^*(\lambda)} P(X^n | A^n, \lambda)}$$

Idea: stop when $GLR_n(\hat{\imath})$ is big for some answer $\hat{\imath}$.

GLR Stopping

For any plausible answer $\hat{\imath} \in i^*(\hat{\mu}(n))$, the GLR_n simplifies to

$$\mathsf{GLR}_n(\hat{\imath}) = \inf_{\lambda:\hat{\imath}\notin i^*(\lambda)} \sum_{a=1}^{K} N_a(n) \,\mathsf{KL}(\hat{\mu}_a(n), \lambda_a)$$

where KL(x, y) is the Kullback-Leibler divergence in the exponential family.

What is a suitable threshold for GLR_n so that we do not make mistakes?

What is a suitable threshold for GLR_n so that we do not make mistakes? A mistake is made when $GLR_n(\hat{\imath})$ is big while $\hat{\imath} \notin i^*(\mu)$.

What is a suitable threshold for GLR_n so that we do not make mistakes? A mistake is made when $GLR_n(\hat{\imath})$ is big while $\hat{\imath} \notin i^*(\mu)$. But then

$$\mathsf{GLR}_n(\hat{\imath}) = \inf_{\lambda:\hat{\imath}\notin i^*(\lambda)} \sum_{a=1}^K N_a(n) \mathsf{KL}(\hat{\mu}_a(n), \lambda_a) \leq \left[\sum_{a=1}^K N_a(n) \mathsf{KL}(\hat{\mu}_a(n), \mu_a)\right]$$

What is a suitable threshold for GLR_n so that we do not make mistakes? A mistake is made when $GLR_n(\hat{\imath})$ is big while $\hat{\imath} \notin i^*(\mu)$. But then

$$\mathsf{GLR}_n(\hat{\imath}) = \inf_{\boldsymbol{\lambda}: \hat{\imath} \notin i^*(\boldsymbol{\lambda})} \sum_{a=1}^K N_a(n) \,\mathsf{KL}(\hat{\mu}_a(n), \lambda_a) \leq \left[\sum_{a=1}^K N_a(n) \,\mathsf{KL}(\hat{\mu}_a(n), \mu_a) \right]$$

Good anytime deviation inequalities exist for that upper bound.

Theorem (Kaufmann and Koolen, 2018)

$$\mathbb{P}\left(\exists n: \sum_{a=1}^{K} N_a(n) \operatorname{KL}(\hat{\mu}_a(n), \mu_a)\right) - \sum_n \ln \ln N_a(n) \ge C(K, \delta)\right) \le \delta$$

for $C(K, \delta) \approx \ln \frac{1}{\delta} + K \ln \ln \frac{1}{\delta}$.

All in all

Final result: lower and upper bound meet on every problem instance.

Theorem (Garivier and Kaufmann 2016)

For the Track-and-Stop algorithm, for any bandit μ

$$\limsup_{\delta \to 0} \frac{\mathbb{E}_{\boldsymbol{\mu}}[\tau]}{\ln \frac{1}{\delta}} = T^*(\boldsymbol{\mu})$$

All in all

Final result: lower and upper bound meet on every problem instance.

Theorem (Garivier and Kaufmann 2016)

For the Track-and-Stop algorithm, for any bandit μ

$$\limsup_{\delta \to 0} \frac{\mathbb{E}_{\boldsymbol{\mu}}[\tau]}{\ln \frac{1}{\delta}} = T^*(\boldsymbol{\mu})$$

Very similar optimality result for *Top Two Thompson Sampling* by Russo (2016). Here $N_i(t)/t \rightarrow w_i^*$ result of posterior sampling.

Problem Variations and Algorithms

Variations

- Prior knowledge about μ
 - Shape constraints: linear, convex, unimodal, etc. bandits
 - Non-parametric (and heavy-tailed) reward distributions (Agrawal, Koolen, and Juneja, 2021)
 - ▶ ...
- Questions beyond Best Arm
 - A/B/n testing (Russac et al., 2021)
 - Robust best arm (part 2 today)
 - Thresholding
 - Best VaR, CVaR and other tail risk measures (Agrawal, Koolen, and Juneja, 2021)
 - ▶ ...
- Multiple correct answers
 - e-best arm
 - In general (Degenne and Koolen, 2019) (Requires a change in lower bound and upper bound)

Lazy Iterative Optimisation of w^*

Instead of computing at every round the plug-in oracle weights

$$w^{*}(\hat{\mu}) = \operatorname{argmax}_{w \in \Delta_{K}} \underbrace{\min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\hat{\mu})} \sum_{i=1}^{K} w_{i} \operatorname{KL}(\hat{\mu}_{i} \| \lambda_{i})}_{\operatorname{concave in } w}$$

We may work as follows

- The inner problem is concave in w.
- It can be maximised iteratively, i.e. with gradient descent.
- We may **interleave sampling** and **gradient** steps.
- A single gradient step per sample is enough (Degenne, Koolen, and Ménard, 2019)

Minimax Action Identification

Model (Teraoka, Hatano, and Takimoto, 20



Maximin Action Identification Problem Find best move at root from samples of leaves.

Model (Teraoka, Hatano, and Takimoto, 20



Maximin Action Identification Problem Find best move at root from samples of leaves.



My Brief History



Best Arm Identification (Garivier and Kaufmann, 2016) Solved, continuous



Depth 2 Game (Garivier, Kaufmann, and Koolen, 2016) Open, continuous?



Depth 1.5 Game (Kaufmann, Koolen, and Garivier, 2018) Solved, discontinuous

What we are able to solve today

Noisy games of any depth



Example Backward Induction Computation



Example Backward Induction Computation



Definition

A game tree is a min-max tree with leaves \mathcal{L} . A bandit model μ assigns a distribution μ_{ℓ} to each leaf $\ell \in \mathcal{L}$.

Definition

A game tree is a min-max tree with leaves \mathcal{L} . A bandit model μ assigns a distribution μ_{ℓ} to each leaf $\ell \in \mathcal{L}$.

The maximin action (best action at the root) is

$$i^*(\mu) := \operatorname{argmax}_{a_1} \min_{a_2} \max_{a_3} \min_{a_4} \cdots \mu_{a_1 a_2 a_3 a_4 \cdots}$$

Definition

A game tree is a min-max tree with leaves \mathcal{L} . A bandit model μ assigns a distribution μ_{ℓ} to each leaf $\ell \in \mathcal{L}$.

The maximin action (best action at the root) is

$$i^*(\mu) := \operatorname{argmax}_{a_1} \min \max_{a_2} \min \cdots_{a_3} \prod_{a_4} \mu_{a_1 a_2 a_3 a_4 \dots}$$

Protocol

For $t = 1, 2, ..., \tau$:

- Learner picks a leaf $L_t \in \mathcal{L}$.
- Learner sees $X_t \sim \mu_{L_t}$

Learner recommends action \hat{I}

Definition

A game tree is a min-max tree with leaves \mathcal{L} . A bandit model μ assigns a distribution μ_{ℓ} to each leaf $\ell \in \mathcal{L}$.

The maximin action (best action at the root) is

$$i^*(\boldsymbol{\mu}) := \operatorname{argmax}_{a_1} \min \max_{a_2} \min \cdots_{a_3} \prod_{a_4} \mu_{a_1 a_2 a_3 a_4 \dots}$$

Protocol

For $t = 1, 2, ..., \tau$:

• Learner picks a leaf $L_t \in \mathcal{L}$.

• Learner sees $X_t \sim \mu_{L_t}$

Learner recommends action \hat{l}

Learner is $\delta\text{-PAC}$ if

$$orall {oldsymbol{\mu}}: \mathop{\mathbb{P}}\limits_{oldsymbol{\mu}} \left(au < \infty \wedge \hat{oldsymbol{l}}
eq i^*(oldsymbol{\mu})
ight) \leq \delta$$

Main Theorem I: Lower Bound

Define the alternatives to μ by Alt $(\mu) = \{\lambda | i^*(\lambda) \neq i^*(\mu)\}$. NB here i^* is best action at the root

Main Theorem I: Lower Bound

Define the alternatives to μ by Alt $(\mu) = \{\lambda | i^*(\lambda) \neq i^*(\mu)\}$. NB here i^* is best action at the root

Theorem (Castro 2014; Garivier and Kaufmann 2016)

Fix a δ -correct strategy. Then for every bandit model μ

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

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

$$\frac{1}{\mathcal{T}^*(\boldsymbol{\mu})} = \max_{\boldsymbol{w} \in \Delta_K} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^K w_i \mathsf{KL}(\boldsymbol{\mu}_i \| \boldsymbol{\lambda}_i).$$

Idea is still to consider the oracle weight map

$$oldsymbol{w}^*(oldsymbol{\mu}) \ \coloneqq \ rgmax_{oldsymbol{w}\in riangle_{K}} \min_{oldsymbol{\lambda}\in riangle \operatorname{Alt}(oldsymbol{\mu})} \ \sum_{i=1}^{K} w_i \operatorname{\mathsf{KL}}(\mu_i \| \lambda_i)$$

and track the plug-in estimate: $L_t \sim w^*(\hat{\mu}(t-1)).$

Idea is still to consider the oracle weight map

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

and track the plug-in estimate: $L_t \sim w^*(\hat{\mu}(t-1)).$

But what about continuity? Does $\hat{\mu}(t) o \mu$ imply $w^*(\mu(t)) o w^*(\mu)$?

Idea is still to consider the oracle weight map

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

and track the plug-in estimate: $L_t \sim w^*(\hat{\mu}(t-1)).$

But what about continuity? Does $\hat{\mu}(t) o \mu$ imply $w^*(\mu(t)) o w^*(\mu)$?

But w^* is not continuous. Even at depth "1.5" with 2 arms.

Idea is still to consider the oracle weight map

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

and track the plug-in estimate: $L_t \sim w^*(\hat{\mu}(t-1)).$

But what about continuity? Does $\hat{\mu}(t)
ightarrow \mu$ imply $w^*(\mu(t))
ightarrow w^*(\mu)$?

But w^* is not continuous. Even at depth "1.5" with 2 arms.

Theorem (Degenne and Koolen, 2019)

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

$$\min_{w \in \boldsymbol{w}^*(\boldsymbol{\mu})} \left\| \boldsymbol{w}_t - \boldsymbol{w} \right\|_\infty \to 0$$

Track-and-Stop is asymptotically optimal.

Example



On Computation

To compute a gradient (in w) we need to differentiate

$$\boldsymbol{w} \mapsto \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^{K} w_i \mathsf{KL}(\boldsymbol{\mu}_i \| \lambda_i)$$

An optimal $\lambda \in Alt(\mu)$ can be found by binary search for common value plus tree reasoning in $O(|\mathcal{L}|)$ (board).

Conclusion

Conclusion

This concludes the guest lecture.

- It has been a pleasure
- Good luck for the exam
- If you have an idea that you want to work on ...

- 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*. Accepted.
- Castro, R. M. (Nov. 2014). "Adaptive sensing performance lower bounds for sparse signal detection and support estimation". In: *Bernoulli* 20.4, pp. 2217–2246.
- Degenne, R. and W. M. Koolen (Dec. 2019). "Pure Exploration with Multiple Correct Answers". In: Advances in Neural Information Processing Systems (NeurIPS) 32, pp. 14591–14600.
- 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, 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. Vol. 2375. Lecture Notes in Computer Science, pp. 255–270.

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)*.

Kaufmann, E. and W. M. Koolen (Oct. 2018). "Mixture Martingales Revisited with Applications to Sequential Tests and Confidence Intervals". Preprint.

 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, pp. 6333–6343.

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. Accepted. Russo, D. (2016). "Simple Bayesian Algorithms for Best Arm Identification". In: CoRR abs/1602.08448.

Teraoka, K., K. Hatano, and E. Takimoto (2014). "Efficient Sampling Method for Monte Carlo Tree Search Problem". In: IEICE Transactions on Infomation and Systems, pp. 392–398.