

# A/B/n Testing with Control in the Presence of Subpopulations

SCOOOL Seminar, INRIA Lille

---

**Wouter M. Koolen**

with Y. Russac, C. Katsimerou, D. Bohle, O. Cappé, A. Garivier

28 Jan 2022



Centrum Wiskunde & Informatica



## Pure Exploration

- We want algorithms for **adaptively** choosing **experiments**
- in an a-priori **unknown environment**
- to gain **information** as fast as possible
- about which decision is **correct/best/useful**



## Pure Exploration

- We want algorithms for **adaptively** choosing **experiments**
- in an a-priori **unknown environment**
- to gain **information** as fast as possible
- about which decision is **correct/best/useful**

⇒ Active Sequential Composite Multiple Hypothesis Test



## Pure Exploration

- We want algorithms for **adaptively** choosing **experiments**
- in an a-priori **unknown environment**
- to gain **information** as fast as possible
- about which decision is **correct/best/useful**

⇒ Active Sequential Composite Multiple Hypothesis Test

## Examples

- $(\epsilon, \delta)$ -PAC learning of policy in unknown MDP
- (MCTS) simulation-based planning
- Drug safety/efficacy trial
- This talk: **A/B/n testing with control**

# The Problem



Control



Version 1



Version 2

...



Version  $K$

# The Problem



Control



Version 1



Version 2

...



Version K



# The Question



Many questions interesting:

BAI	What is the <b>best</b> version?	$\{0, \dots, K\}$
Thr.	Which versions are at least $\theta$ -good?	$\mathcal{P}(\{0, \dots, K\})$
	• Is <b>any</b> version better than the control?	$\{\text{yes, no}\}$
	• <b>Which</b> version, if any, beats the control?	$\{0, \dots, K\}$
ABC	<b>Which versions</b> are better than the control?	$\mathcal{P}(\{1, \dots, K\})$

this  
talk

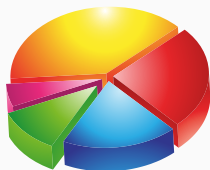
How does the presence of **sub-populations** affect learning?

# The Model

## Definition (Model)

A bandit with  $K + 1$  arms and  $J$  subpopulations is

- Distribution  $\alpha$  of the  $J$  subpopulations
- a  $(K + 1) \times J$  matrix  $\mu$  of (Bernoulli, say) reward distributions



$\alpha$



$\mu$



# The Target

## All arms Better than the Control with Subpopulations (ABC-S)

The **correct answer** for Bandit  $(\alpha, \mu)$  is

$$S^*(\mu) = \left\{ k \in \{1, \dots, K\} \mid \sum_{j=1}^J \alpha_j \mu_{k,j} > \sum_{j=1}^J \alpha_j \mu_{0,j} \right\}$$

# The Protocol

We study four **Modes of Interaction**

## Protocol

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

- |                                 |                       |                       |                      |
|---------------------------------|-----------------------|-----------------------|----------------------|
| <b>Oblivious</b>                | <b>Agnostic</b>       | <b>Propert.</b>       | <b>Active</b>        |
| Pick $A_t$                      | Pick $A_t$            | See $I_t \sim \alpha$ |                      |
| <b>Hidden</b> $I_t \sim \alpha$ | See $I_t \sim \alpha$ | Pick $A_t$            | Pick $A_t$ and $I_t$ |
- See reward  $X_t \sim \mu_{A_t, I_t}$

Learner recommends  $\hat{S} \subseteq \{1, \dots, K\}$  (arms better than control)

Modes *constrain* the sampling proportions of  $(I, A)$



We want our learner to

(1) be  $\delta$ -PAC, i.e. for any bandit  $\mu$ ,

$$\mathbb{P}_{\mu}(\text{Learner stops and recommends wrong answer}) \leq \delta.$$

(2) minimise **sample complexity**, i.e.  $\mathbb{E}_{\mu}$  [stopping time]

# Our Results

(Russac, Katsimerou, Bohle, Cappé, Garivier, and Koolen, 2021)

- Information-theoretic lower bounds for all four modes
- **Matching** ( $\delta \rightarrow 0$ ) algorithms (Track-and-Stop family)

# Results

Let's think about  $K = 1$  arm vs control with  $J$  subpopulations.

Let's investigate the **Gaussian** case: reward for  $a, j$  is

$$\mathcal{N}(\mu_{a,j}, \sigma_{a,j}^2).$$

# Lower Bound

## Theorem

For any strategy, the expected number of rounds for the ABC-S problem with mode constraint  $\mathcal{C}$  satisfies

$$\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau_{\delta}]}{\ln(1/\delta)} \geq T^*(\mu)$$

where

$$\begin{aligned} T^*(\mu)^{-1} &= \max_{w \in \mathcal{C}} \inf_{\lambda: S^*(\lambda) \neq S^*(\mu)} \sum_{a=0}^K \sum_{i=1}^J w_{a,i} \text{KL}(\mu_{a,i}, \lambda_{a,i}) \\ &= \max_{w \in \mathcal{C}} \min_{b \neq 0} \inf_{\substack{\lambda \in \mathcal{L} \\ \alpha^T \lambda_0 = \alpha^T \lambda_b}} \sum_{a \in \{0, b\}} \sum_{i=1}^J w_{a,i} \text{KL}(\mu_{a,i}, \lambda_{a,i}) \end{aligned}$$

NB: the min/inf is the (expected) amount of statistical evidence collected per round by sampling proportions  $w$  against any bandit  $\lambda$  with  $S^*(\lambda) \neq S^*(\mu)$

## How to think about this

Estimate for arm quality carries **uncertainty**:

$$\sum_{j=1}^J \alpha_j \hat{\mu}_{j,a}$$

Uncertainty  $\Leftrightarrow$  **variance**. Now if  $(a, j)$  is sampled  $n_{a,j}$  times,

$$\mathbb{V} \left[ \sum_{j=1}^J \alpha_j \hat{\mu}_{j,a} \right] = \sum_{j=1}^J \alpha_j^2 \mathbb{V} [\hat{\mu}_{j,a}] = \sum_{j=1}^J \frac{\alpha_j^2 \sigma_{a,j}^2}{n_{a,j}}$$

## How to think about this

Estimate for arm quality carries **uncertainty**:

$$\sum_{j=1}^J \alpha_j \hat{\mu}_{j,a}$$

Uncertainty  $\Leftrightarrow$  **variance**. Now if  $(a, j)$  is sampled  $n_{a,j}$  times,

$$\mathbb{V} \left[ \sum_{j=1}^J \alpha_j \hat{\mu}_{j,a} \right] = \sum_{j=1}^J \alpha_j^2 \mathbb{V} [\hat{\mu}_{j,a}] = \sum_{j=1}^J \frac{\alpha_j^2 \sigma_{a,j}^2}{n_{a,j}}$$

Minimised unconstrained (active mode) at

$$n_{a,j} \propto \alpha_j \sigma_{a,j}$$

Other modes: add constraints  $\mathbf{n} \in \mathcal{C}$



## Results (explicit Gaussian case)

$K = 1$  arm vs control with  $J$  subpopulations.

Denoting the gap by  $\Delta_1 = \sum_{j=1}^J \alpha_j (\mu_{1,j} - \mu_{0,j})$ , we find

$$T_{\text{oblivious}}^*(\mu) \approx \frac{2 \left( \sum_{a \in \{0,1\}} \sqrt{\sum_{j=1}^J \alpha_j (\sigma_{a,j}^2 + (\mu_{a,j} - \mu_a)^2)} \right)^2}{\Delta_1^2}$$

$$T_{\text{agnostic}}^*(\mu) = \frac{2 \left( \sqrt{\sum_{j=1}^J \alpha_j \sigma_{0,j}^2} + \sqrt{\sum_{j=1}^J \alpha_j \sigma_{1,j}^2} \right)^2}{\Delta_1^2}$$

$$T_{\text{proport.}}^*(\mu) = \frac{2 \sum_{j=1}^J \alpha_j (\sigma_{0,j} + \sigma_{1,j})^2}{\Delta_1^2},$$

$$T_{\text{active}}^*(\mu) = \frac{2 \left( \sum_{j=1}^J \alpha_j (\sigma_{0,j} + \sigma_{1,j}) \right)^2}{\Delta_1^2},$$

better



# Algorithm

## Sampling Rule

Ensure that actual sampling proportions  $\mathbf{N}_t/t$  track oracle proportions at **plug-in estimate**  $\hat{\mu}(t)$

$$w^*(\hat{\mu}(t)) = \arg \max_{w \in \mathcal{C}} \min_{b \neq 0} \inf_{\substack{\lambda \in \mathcal{L} \\ \alpha^\top \lambda_0 = \alpha^\top \lambda_b}} \sum_{a \in \{0, b\}} \sum_{i=1}^J w_{a,i} \text{KL}(\hat{\mu}_{a,i}(t), \lambda_{a,i})$$

Tracking is done locally, respecting the mode *constraint*

## Stopping Rule

Stop at  $\tau_\delta = t$  when we've collected enough information, i.e.

$$\min_{b \neq 0} \inf_{\substack{\lambda \in \mathcal{L} \\ \alpha^\top \lambda_0 = \alpha^\top \lambda_b}} \sum_{a \in \{0, b\}} \sum_{i=1}^J N_{a,i}(t) \text{KL}(\hat{\mu}_{a,i}(t), \lambda_{a,i}) \geq \ln \frac{\ln t}{\delta}$$

## Recommendation Rule

Output  $S^*(\hat{\mu}(t))$

# Validation: Asymptotic Optimality

## Theorem

*The stopping+recommendation rules are  $\delta$ -PAC.*

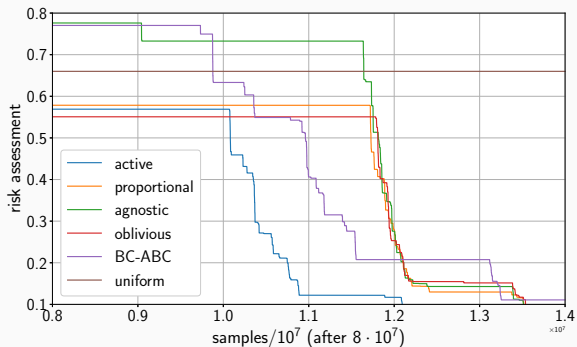
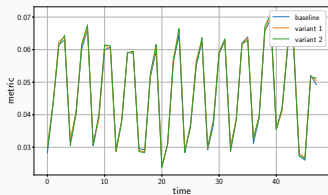
## Theorem

*The algorithm ensures that the expected number of rounds for the ABC-S problem with mode constraint  $\mathcal{C}$  satisfies*

$$\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau_{\delta}]}{\ln(1/\delta)} \leq T^*(\mu)$$

Upper bound matching lower bound, perfectly.

# Validation: Real-world experiment with $\delta = 0.1$



# Conclusion

- Interesting pure exploration problems in A/B/n testing
- Subpopulation **awareness** reduces **sample complexity** ...  
... even if the recommendations (arms) are **unaware**!

# Conclusion

- Interesting pure exploration problems in A/B/n testing
- Subpopulation **awareness** reduces **sample complexity** ...  
... even if the recommendations (arms) are **unaware**!



Next steps:

- Going beyond asymptotic optimality
- Structured (shape-constrained) mean matrices

# Thanks!

## References

---

-  Garivier, A. and E. Kaufmann (2016). “Optimal best arm identification with fixed confidence”. In: *Conference on Learning Theory*. PMLR, pp. 998–1027.
-  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.