

# Choosing Answers in $\varepsilon$ -Best-Answer Identification for Linear Bandits

Marc Jourdan and Rémy Degenne

June 17, 2022



**Initial goal:** Identify the item having the highest averaged return.

**Problem:** When the two best items have highly similar averaged return, the number of samples required to differentiate them is large.

**Corrected goal:** Identify one item which is  $\varepsilon$ -close to the best one ( $\varepsilon$ -BAI).

**Challenge:** Multiple correct answers.

? How to choose among the set of  $\varepsilon$ -optimal answers ?

👉 Focus on the  $\varepsilon$ -optimal answer which is the easiest to verify.

Transductive linear Gaussian bandits:

- arm  $a \in \mathcal{K}$ , finite subset of  $\mathbb{R}^d$ ,
- answer  $z \in \mathcal{Z}$ , finite subset of  $\mathbb{R}^d$ ,
- unknown mean parameter,  $\mu \in \mathbb{R}^d$ .

At time  $t$ , pull  $a_t \in \mathcal{K}$  and observe  $X_t^{a_t} \sim \mathcal{N}(\langle \mu, a_t \rangle, 1)$ .

**Goal:** Identify one  $\varepsilon$ -optimal answer with confidence  $\delta$ ,  $z \in \mathcal{Z}_\varepsilon(\mu)$ .

**Objective:** Minimize  $\mathbb{E}_\mu[\tau_\delta]$  for  $(\varepsilon, \delta)$ -PAC algorithms

$$\mathbb{P}_\mu[\tau_\delta < +\infty, z_{\tau_\delta} \notin \mathcal{Z}_\varepsilon(\mu)] \leq \delta.$$

? What is the best one could achieve? Degenne and Koolen (2019)

👉 For all  $(\varepsilon, \delta)$ -PAC strategy, for all  $\mu$ ,  $\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_\mu[\tau_\delta]}{\ln(1/\delta)} \geq T_\varepsilon(\mu)$ .

? How to choose among the set of  $\varepsilon$ -optimal answers ?

👉 **Furthest answer**:  $\varepsilon$ -optimal answer for which its alternative is the easiest to differentiate from thanks to an optimal allocation over arms.

$$(z_F(\mu), w_F(\mu)) \stackrel{\text{def}}{=} \arg \max_{(z,w) \in \mathcal{Z}_\varepsilon(\mu) \times \Delta_K} \inf_{\lambda \in \neg_\varepsilon z} \frac{1}{2} \|\mu - \lambda\|_{V_w}^2,$$

are the maximizers realizing  $T_\varepsilon(\mu)$ .  $\neg_\varepsilon z$  alternative to  $z$ ,  $\Delta_K$  simplex,  $V_w = \sum_{a \in \mathcal{K}} w^a a a^\top$  design matrix with norm  $\|\cdot\|_{V_w}$ .

**Greedy answer**:  $z^*(\mu) = \arg \max_{z \in \mathcal{Z}} \langle \mu, z \rangle$ , unique correct answer in BAI.

👉 sample inefficient, 10% higher empirical stopping time for  $\delta = 1\%$ .

# Adapting any BAI algorithm for $\varepsilon$ -BAI

? How to stop to obtain an  $(\varepsilon, \delta)$ -PAC strategy ?

👉 calibrated **GLR stopping rule** for  $z_t \in \mathcal{Z}_\varepsilon(\mu_{t-1})$

? Which  $z_t$  should we **recommend** to stop as early as possible ?

👉 **Instantaneous furthest answer**:  $\varepsilon$ -optimal answer with highest GLR

$$z_F(\mu_{t-1}, N_{t-1}) = \arg \max_{z \in \mathcal{Z}_\varepsilon(\mu_{t-1})} \inf_{\lambda \in \neg_\varepsilon z_t} \|\mu_{t-1} - \lambda\|_{V_{N_{t-1}}}^2,$$

where  $N_{t-1}^a = \sum_{s=1}^{t-1} \mathbf{1}_{\{a_s=a\}}$  and  $\mu_{t-1} = V_{N_{t-1}}^{-1} \sum_{s=1}^{t-1} X_s^{a_s} a_s$ .

? How to **modify any BAI algorithms** to be  $(\varepsilon, \delta)$ -PAC ?

👉 use GLR stopping with  $z_t \in z_F(\mu_{t-1}, N_{t-1})$ ,

👉 keep the sampling rule unchanged.

# $L_\epsilon$ BAI (Linear $\epsilon$ -BAI)

Can we achieve asymptotic optimality and be empirically competitive ?

👉  $L_\epsilon$ BAI, by using the concept of furthest answer in the sampling rule.

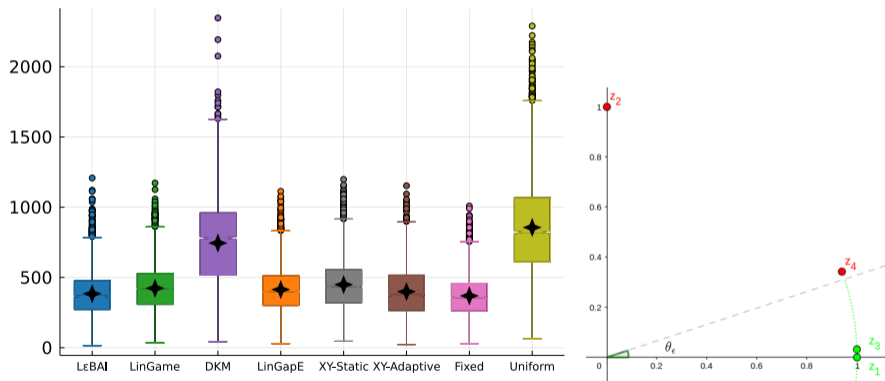


Figure: Empirical stopping time of  $L_\epsilon$ BAI compared to modified BAI algorithms.

# Conclusion

- 1 Don't choose greedily: aim at identifying the *furthest* answer !
- 2 Simple procedure to adapt your favorite BAI algorithm to  $\varepsilon$ -BAI.
- 3  $L\varepsilon$ BAI, asymptotically optimal and empirically competitive.

Paper & Poster

