

Anytime Model Selection in Linear Bandits

Parnian Kassraie, Aldo Pacchiano, Nicolas Emmenegger, Andreas Krause



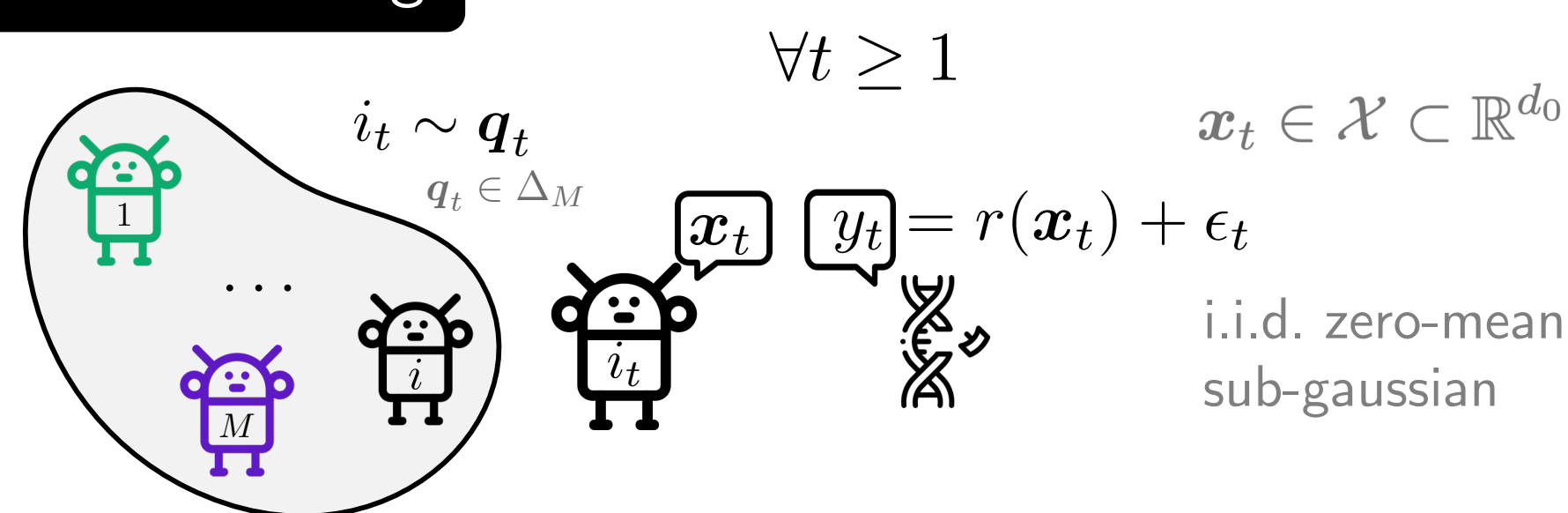
ETH zürich

Can we perform adaptive model selection, while simultaneously optimizing for a reward? Can we be sample-efficient & anytime?

- Solving a Bandit/BO problem:
 - Commit to a reward model (a priori)
 - Interact with the environment accordingly
- There are many ways to model the reward

$$M \gg n \quad n: \text{horizon/stopping time}$$
- Not known a priori which agent is going to be the best (e.g. in terms of sample efficiency, or regret)
- We can select the model based on empirical evidence.

Problem Setting



- Model Class

$$\{\phi_j : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^d, j = 1, \dots, M\} \quad + \text{typical regularity assumptions}$$

$$\exists j^* \in [M] \text{ s.t. } r(\cdot) = \theta_{j^*}^\top \phi_{j^*}(\cdot)$$
- Agents

Agent j only uses ϕ_j to model the reward

Update its action selection policy $p_{t,j} \in \mathcal{M}(\mathcal{X})$

Using the full history $H_{t-1} = \{(x_1, y_1), \dots, (x_{t-1}, y_{t-1})\}$
- Goal

$$R(n) = \sum_{t=1}^n r(x^*) - r(x_t) \quad n \text{ unknown}$$

Ingredient I: Exponential Weights Updates

Adjust the probability of selecting each agent wrt the reward it has obtained so far

- known to yield $\log M$ regret in full-info setting

$$q_{t,j} = \frac{\exp(\eta_t \sum_{s=1}^{t-1} \hat{r}_{s,j})}{\sum_{i=1}^M \exp(\eta_t \sum_{s=1}^{t-1} \hat{r}_{s,i})}$$

η_t sensitivity of updates

- we don't observe the reward of every agent, so:

$$\hat{r}_{t,j} = \mathbb{E}_{x \sim p_{t,j}} [\hat{\theta}_t^\top \phi(x)]$$

- Regret will depend on bias and variance of $\hat{\theta}_t$
- Typical online regression oracles are $\sqrt{M} \rightarrow \text{poly} M$ regret

Ingredient II: Sparse Online Regression Oracle

Turn lasso into a **sparse** online regression oracle

$$\hat{\theta}_t = \arg \min_{\theta} \frac{1}{t} \|\mathbf{y}_t - \Phi_t \theta\|_2^2 + \lambda_t \sum_{j=1}^M \|\theta_j\|_2$$

Bias and variance are both $\log M$

Theorem (Anytime Conf. Seq.)

If for all $t \geq 1$

$$\lambda_t \geq \frac{c_1}{\sqrt{t}} \sqrt{\log(M/\delta) + \sqrt{d} (\log(M/\delta) + (\log \log d)_+)}$$

cost of going 'time uniform'

then,

$$\mathbb{P} \left(\forall t \geq 1 : \|\theta - \hat{\theta}_t\|_2 \leq \frac{c_2 \lambda_t}{\kappa^2(\Phi_t, 2)} \right) \geq 1 - \delta$$

c_1 and c_2 made exact in the paper

Restricted Eigenvalue property [check paper]

Simultaneous Model Selection and Optimization

- Putting it together we get **Anytime Exponential** weighting algorithm with Lasso reward estimates (ALEXP)

Algorithm 1 ALEXP

Inputs: $\gamma_t, \eta_t, \lambda_t$ for $t \geq 1$

for $t \geq 1$ **do**

Draw $\mathbf{x}_t \sim (1 - \gamma_t) \sum_{j=1}^M q_{t,j} p_{t,j} + \gamma_t \text{Unif}(\mathcal{X})$ mix with exploration

Observe $y_t = r(\mathbf{x}_t) + \epsilon_t$.

Append history $H_t = H_{t-1} \cup \{(\mathbf{x}_t, y_t)\}$.

Update agents $p_{t,j}$ for $j = 1, \dots, M$.

Calculate $\hat{\theta}_t \leftarrow \text{Lasso}(H_t, \lambda_t)$ and estimate

$$\hat{r}_{t,j} \leftarrow \mathbb{E}_{x \sim p_{t+1,j}} [\hat{\theta}_t^\top \phi(x)]$$

Update selection distribution



$$q_{t+1,j} \leftarrow \frac{\exp(\eta_t \sum_{s=1}^t \hat{r}_{s,j})}{\sum_{i=1}^M \exp(\eta_t \sum_{s=1}^t \hat{r}_{s,i})}$$

Theorem (Regret - Informal)

For appropriate choices of $(\gamma_t, \lambda_t, \eta_t)$, prescribed in the paper

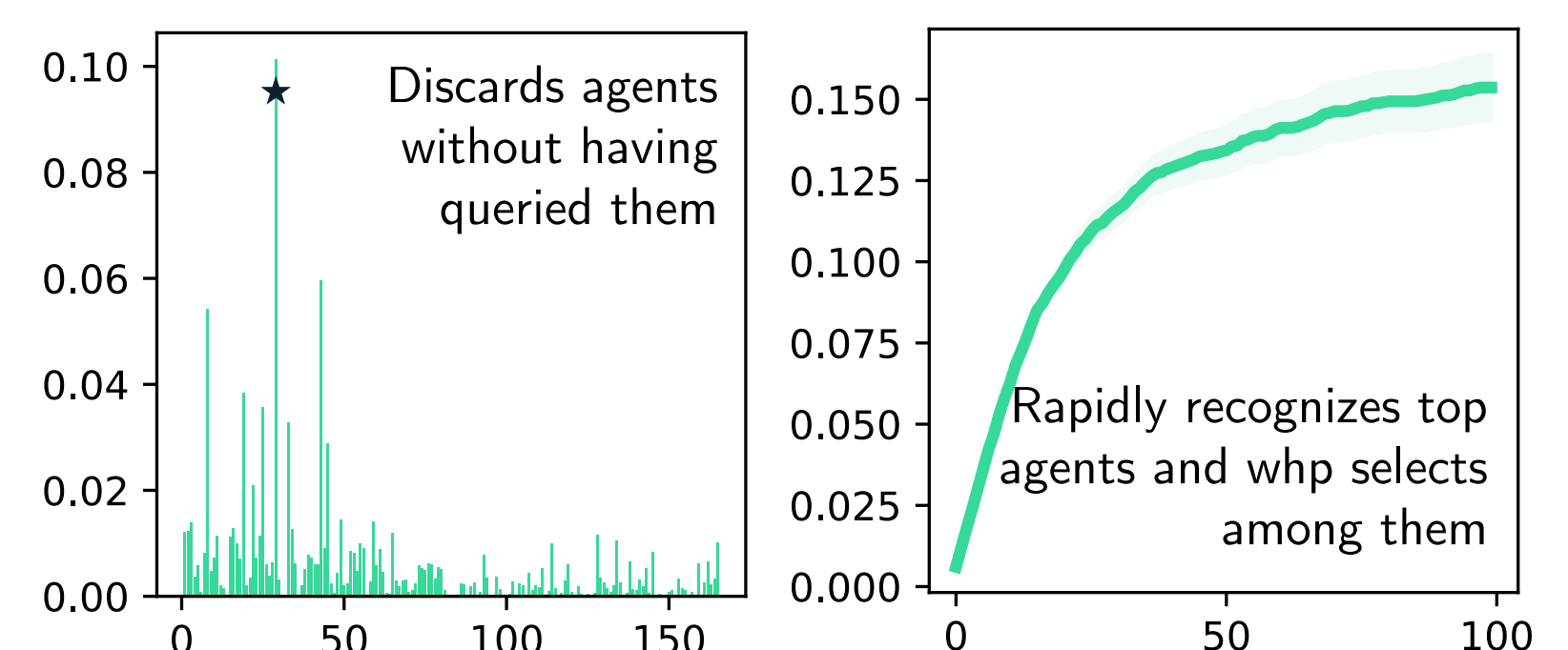


$$R(n) = \mathcal{O} \left(C(M, \delta, d) \left(\sqrt{n} \log M + n^{3/4} \right) \right)$$

with probability greater than $1 - \delta$, simultaneously for all $n \geq 1$.

$$C(M, \delta, d) = \mathcal{O} \left(\sqrt{d \log M / \delta} + \sqrt{d \log M / \delta} \right)$$

Model Selection Dynamics



Comparison to prior work

| | technique | $\log M$ regret | MS guarantee | adaptive & anytime |
|------------------------------|----------------------|-----------------|--------------|--------------------|
| Sparse Linear Bandits | Lasso | ✓ | ✗ | ✗ |
| MS for Black-Box Bandits | OMD with bandit info | ✗ | ✓ | ✗ |
| MS for Linear Bandits (Ours) | EXP4 with full info | ✓ | ✓ | ✓ |

