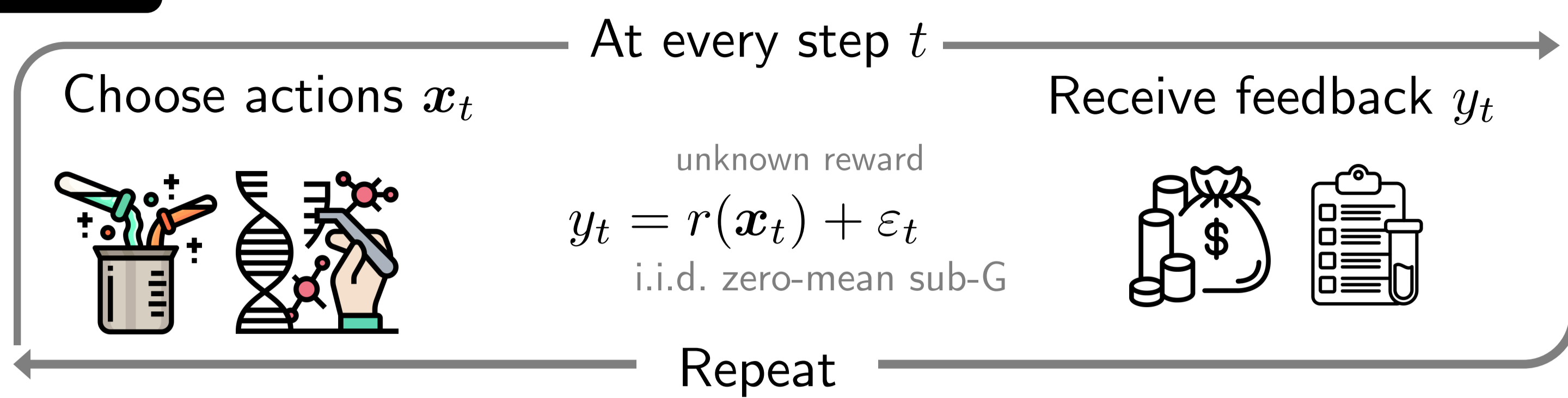


Anytime Model Selection in Linear Bandits

Parnian Kassaie, Nicolas Emmenegger, Andreas Krause, Aldo Pacchiano



Intro



- The statistical modeling of the reward function plays a crucial role in efficiency of bandit algorithms -- they maintain an estimate of the target function, and use it to choose the next action.

- It is not known a priori which model is going to yield the most sample efficient algorithm, and we can only select the right model as we gather empirical evidence.

- Online Model Selection is not fun and games. $\mathbf{x}_t \in \mathcal{X} \subset \mathbb{R}^{d_0}$

$$H_{t-1} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1})\}$$

Reward maximization → not so diverse sample

History dependence → non-i.i.d sample

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

- Our setting $\{\phi_j : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^d, j = 1, \dots, M\}$

$$\exists j^* \in [M] \text{ s.t. } r(\cdot) = \boldsymbol{\theta}_{j^*}^\top \phi_{j^*}(\cdot)$$

$$M \gg T$$

+ typical regularity assumptions

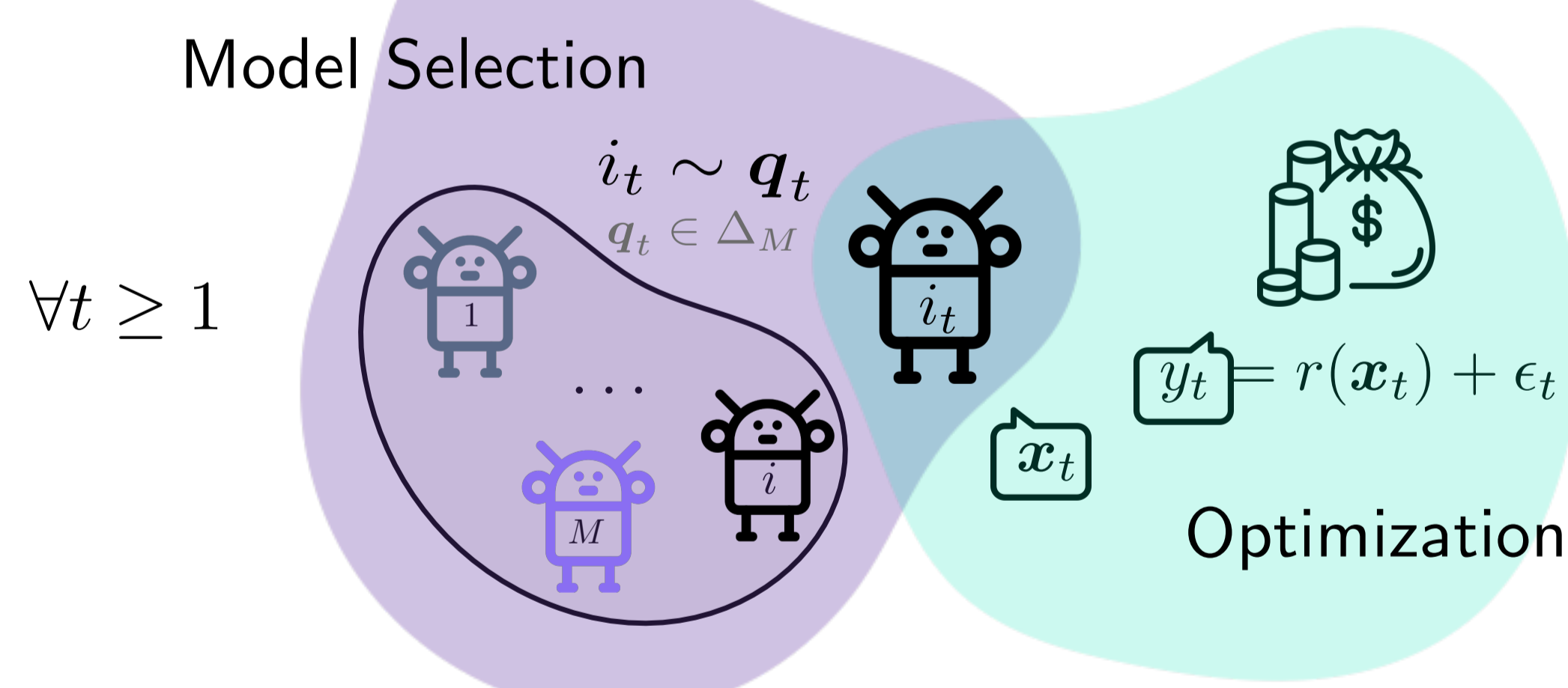
- Online Model Selection problem

Find j^* while maximizing for the unknown r

$$R(T) = \sum_{t=1}^T r(\mathbf{x}^*) - r(\mathbf{x}_t) \quad \begin{array}{l} \text{-- Sublinear in } T \\ \text{-- } \log M \end{array}$$

Approach

- Probabilistic Aggregation: Instantiate M algorithm each using a different ϕ_j to model the reward. Randomly iterate over them.



- With probability $q_{t,j}$ choose agent j and let them choose an action according to their action selection policy $p_{t,j} \in \mathcal{M}(\mathcal{X})$

Ingredient I: Exponential Weights Updates

- Increase $q_{t,j}$ if the the agent *seems* to be lucrative

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

Estimate of the reward obtained by agent j so far

sensitivity of updates

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

- This technique is known to yield $\log M$ regret in full-info setting, when all $r_{t,j}$ are known. But now, the regret will depend on the bias and variance of $\hat{\boldsymbol{\theta}}_t$

- Typical online regression oracles are $\sqrt{M} \rightarrow \text{poly}M$ regret

Main Results

- This gives ALEXP: Anytime Exponential weighting algorithm with Lasso reward estimates

log M regret	MS guarantee	adaptive & anytime
✓	✓	✓

Algorithm 1 ALEXP

Inputs: $\gamma_t, \eta_t, \lambda_t$ for $t \geq 1$ **for** $t \geq 1$ **do**Draw $\mathbf{x}_t \sim (1 - \gamma_t) \sum_{j=1}^M q_{t,j} p_{t,j} + \gamma_t \text{Unif}(\mathcal{X})$ Observe $y_t = r(\mathbf{x}_t) + \varepsilon_t$.Append history $H_t = H_{t-1} \cup \{(\mathbf{x}_t, y_t)\}$.Update agents $p_{t,j}$ for $j = 1, \dots, M$.Calculate $\hat{\boldsymbol{\theta}}_t \leftarrow \text{Lasso}(H_t, \lambda_t)$ and estimate $\hat{r}_{t,j}$ Update selection distribution $q_{t+1,j}$ **end for**

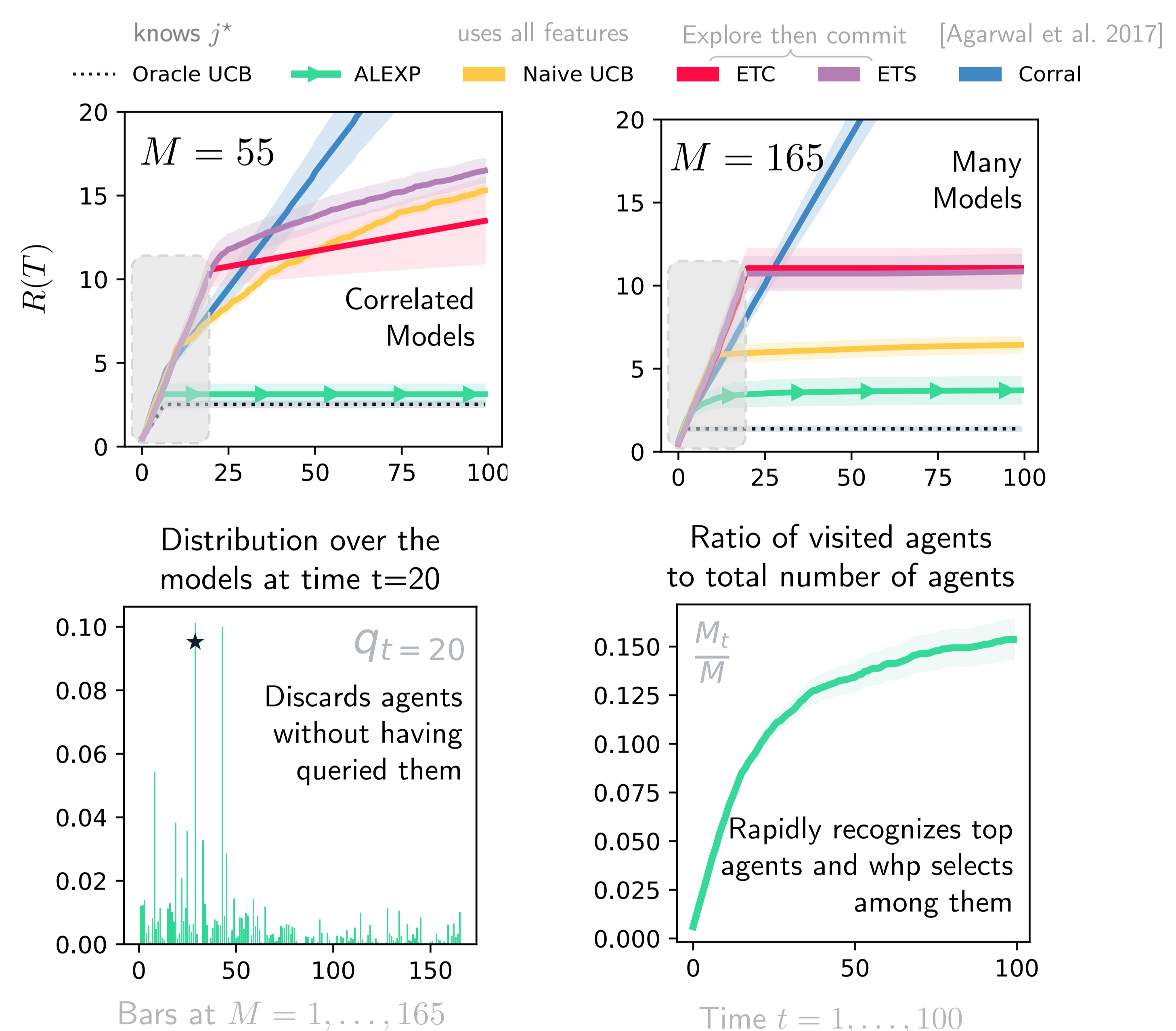
Theorem (Regret - Informal)

For appropriate choices of parameters,
[prescribed in the paper]

$$R(T) = \tilde{O}\left(\sqrt{T \log^3 M} + T^{3/4} \sqrt{\log M}\right)$$

w.h.p. simultaneously for all $T \geq 1$.

Empirical Insights



Ingredient II: Sparse Online Regression Oracle

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

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

Theorem (Anytime Lasso Conf Seq)

If for all $t \geq 1$ Bias and variance are both $\log M$

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

cost of going 'time uniform'

then,

Restricted Eigenvalue property

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