



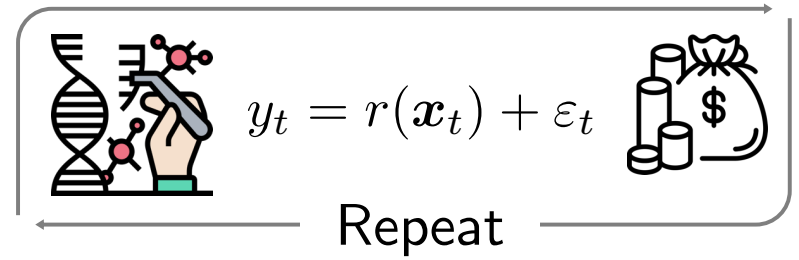
Anytime Model Selection for Linear Bandits

Parnian Kassraie, Nicolas Emmenegger, Andreas Krause, Aldo Pacchiano



Anytime Model Selection

At every step t

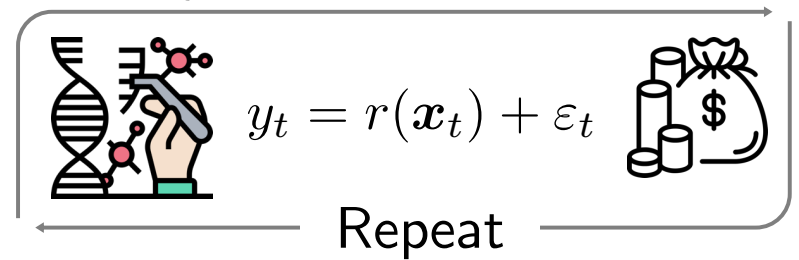


Anytime Model Selection

Solving a Linear Bandit problem :

1. Commit to a reward model (a priori)
2. Interact with the environment to maximize reward

At every step t



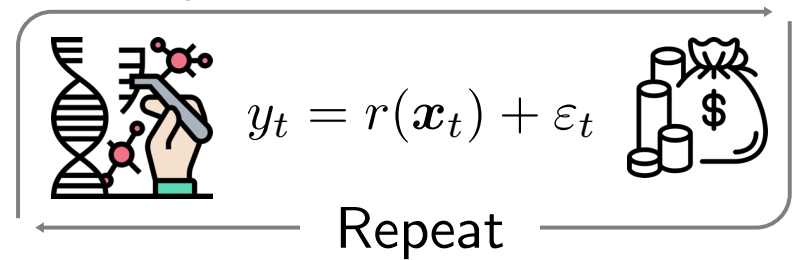
Anytime Model Selection

Solving a Linear Bandit problem :

1. Commit to a reward model (a priori)
2. Interact with the environment to maximize reward

There are many ways to model r

At every step t



Anytime Model Selection

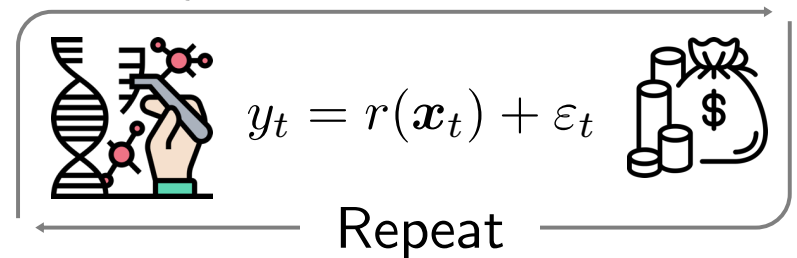
Solving a Linear Bandit problem :

1. Commit to a reward model (a priori)
2. Interact with the environment to maximize reward

There are **many** ways to model r

$$\{\phi_j : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^d, j = 1, \dots, M\}$$
$$\exists j^* \in [M] \text{ s.t. } r(\cdot) = \boldsymbol{\theta}_{j^*}^\top \phi_{j^*}(\cdot)$$

At every step t



$$M \gg T \text{ horizon/stopping time}$$

Anytime Model Selection

Solving a Linear Bandit problem :

1. Commit to a reward model (a priori)
2. Interact with the environment to maximize reward

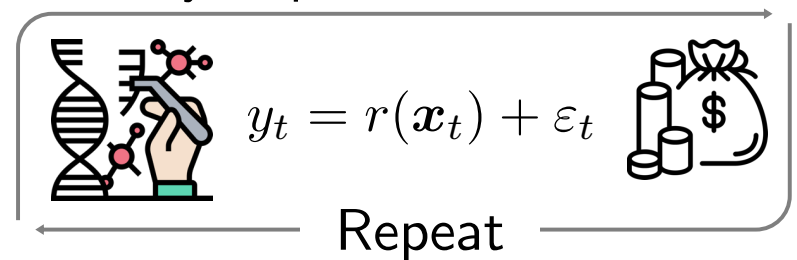
There are many ways to model r

$$\{\phi_j : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^d, j = 1, \dots, M\}$$
$$\exists j^* \in [M] \text{ s.t. } r(\cdot) = \theta_{j^*}^\top \phi_{j^*}(\cdot)$$

Not known a priori which model is going to yield the best algo.

... but we can guess based on empirical evidence.

At every step t



$$M \gg T \text{ horizon/stopping time}$$

Anytime Model Selection

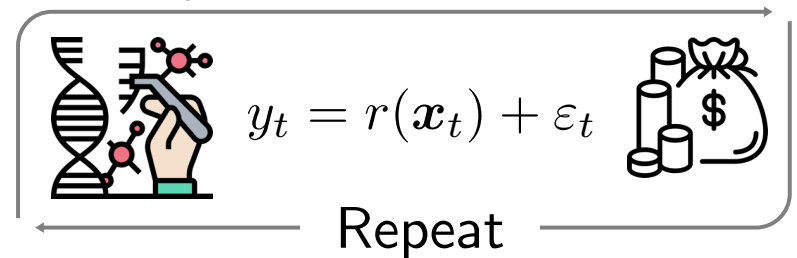
Solving a Linear Bandit problem :

1. Commit to a reward model (a priori)
2. Interact with the environment to maximize reward

There are many ways to model r

$$\{\phi_j : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^d, j = 1, \dots, M\}$$
$$\exists j^* \in [M] \text{ s.t. } r(\cdot) = \theta_{j^*}^\top \phi_{j^*}(\cdot)$$

At every step t



$M \gg T$ horizon/stopping time

Not known a priori which model is going to yield the best algo.
... but we can guess based on empirical evidence.

Anytime Model Selection problem

Find j^* while maximizing for the unknown r

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

Online Model Selection problem

Find j^* while maximizing for the unknown r

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

Online Model Selection problem

Find j^*

$\forall T \geq 1$

$t=1$

$- \log M$

linear in T

r

Why do we need to select?
Why not just try out everything?

Online Model Selection problem

Find j^*

$\forall T \geq 1$

$t=1$

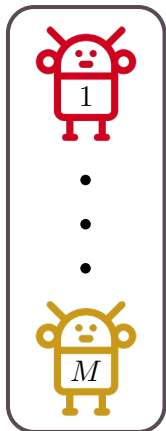
$- \log M$

linear in T

r

Why do we need to select?
Why not just try out everything?

Instantiate M algorithms each using a different model



Online Model Selection problem

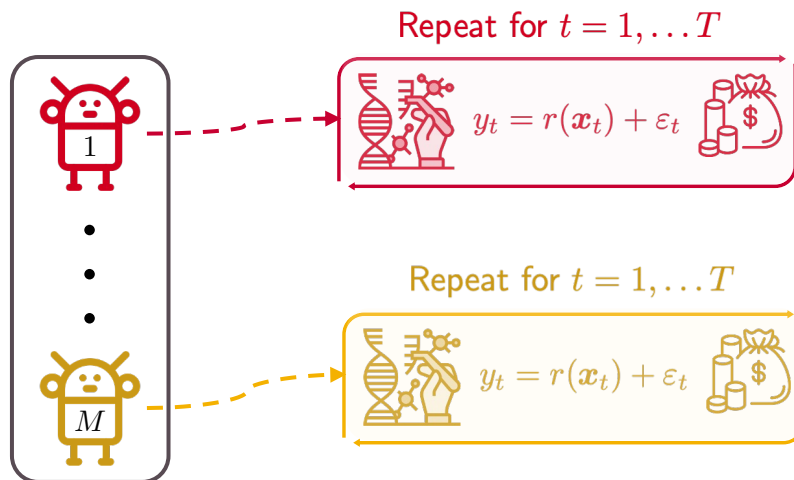
$$\text{Find } j^* \text{ such that } \sum_{t=1}^T r_{j_t} - \min_{j \in M} \sum_{t=1}^T r_{j_t} \leq \epsilon$$

$\forall T \geq 1$

Why do we need to select?
Why not just try out everything?

Instantiate M algorithms each using a different model

Run **all** algorithms in parallel



Online Model Selection problem

$$\text{Find } j^* \\ \forall T \geq 1$$

Why do we need to select?
Why not just try out everything?

$$t=1$$

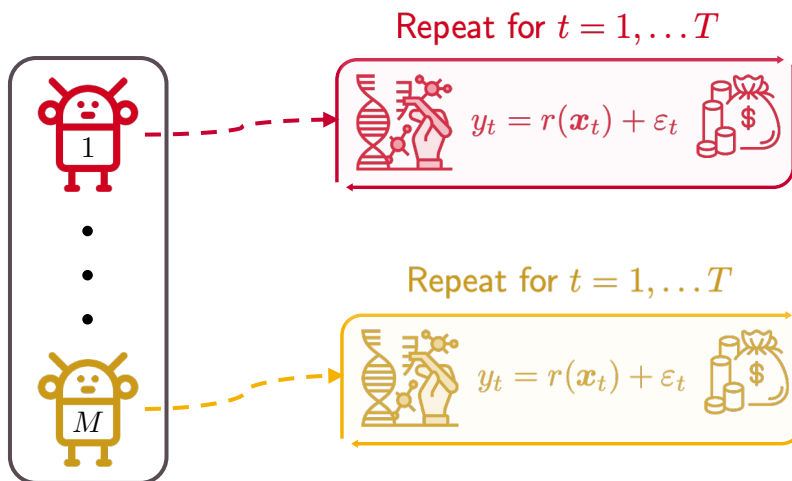
$$- \log M$$

$n r$

linear in T

Instantiate M algorithms each using a different model

Run **all** algorithms in parallel



Statistically expensive
 \leftrightarrow High regret

$$\text{poly}(M)$$

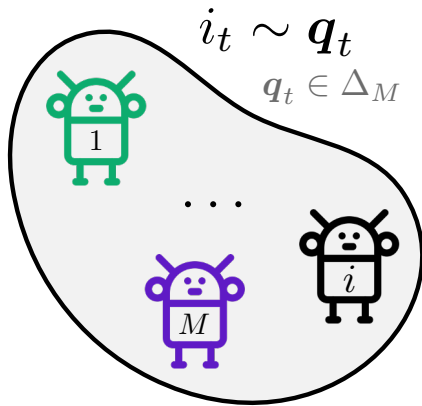
w.h.p.

Our Solution: Probabilistic Aggregation

 Randomly iterate over the agents and at each step play only one

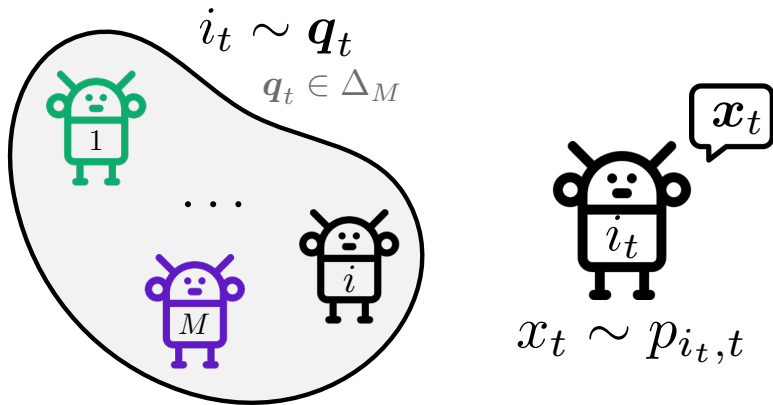
Our Solution: Probabilistic Aggregation

💡 Randomly iterate over the agents and at each step play only one



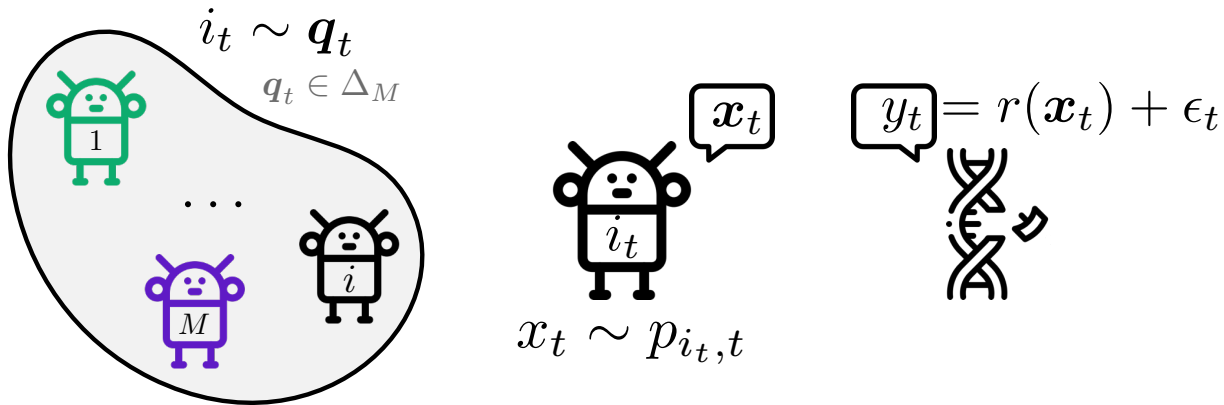
Our Solution: Probabilistic Aggregation

💡 Randomly iterate over the agents and at each step play only one



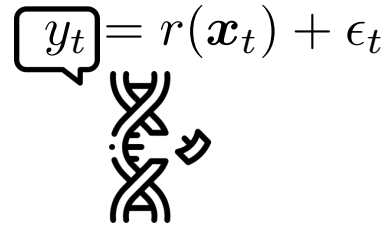
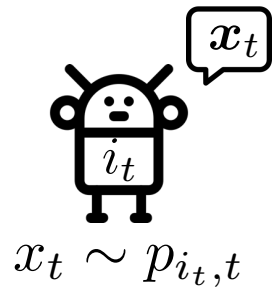
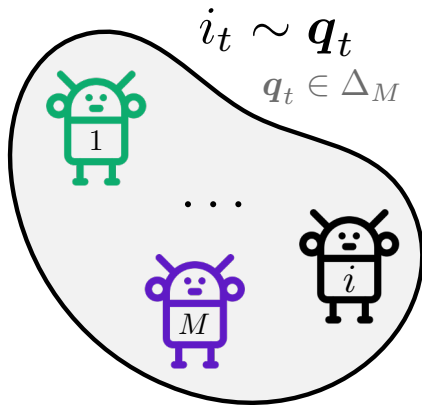
Our Solution: Probabilistic Aggregation

💡 Randomly iterate over the agents and at each step play only one



Our Solution: Probabilistic Aggregation

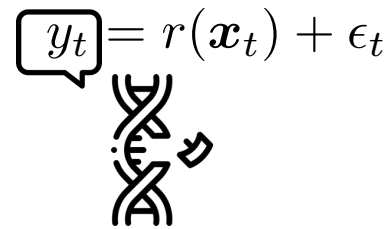
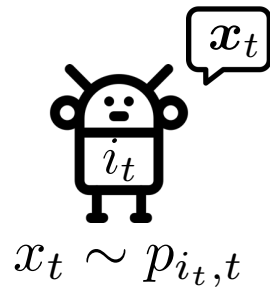
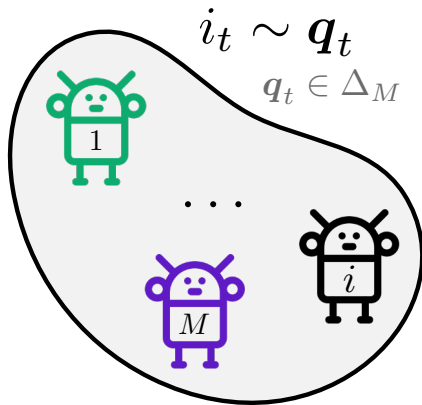
💡 Randomly iterate over the agents and at each step play only one



Update all agents
Update \mathbf{q}_t

Our Solution: Probabilistic Aggregation

💡 Randomly iterate over the agents and at each step play only one



Update all agents
Update \mathbf{q}_t

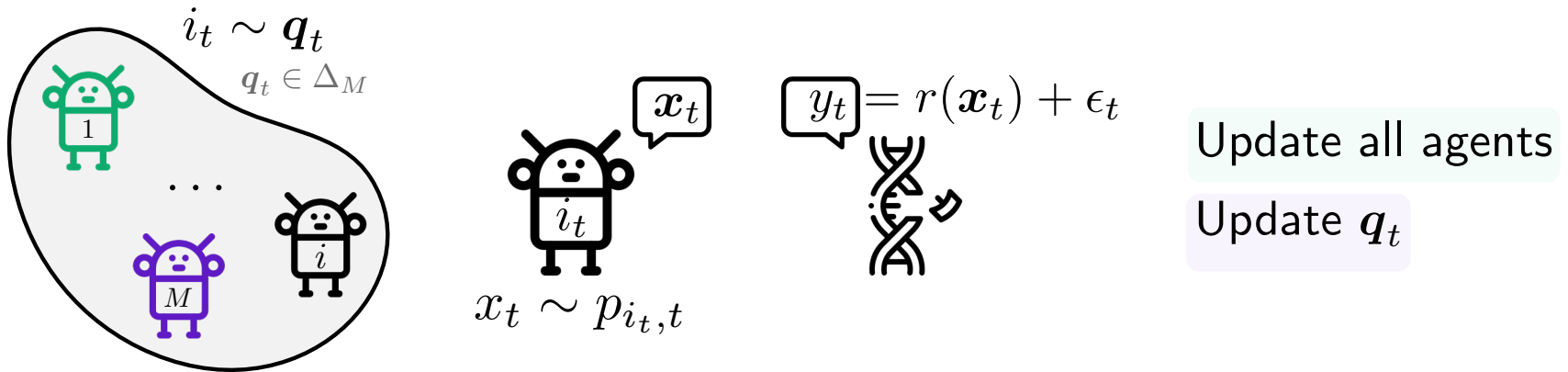
Play one agent, but update all.
Reward not observed? **Estimate** it.

Choose your **estimator** very carefully!

$$\hat{r}_{t,j} \text{ for } j = 1, \dots, M$$

Our Solution: Probabilistic Aggregation

💡 Randomly iterate over the agents and at each step play only one



Play one agent, but update all.
Reward not observed? **Estimate** it.

$$\hat{r}_{t,j} \text{ for } j = 1, \dots, M$$

Choose your **estimator** very carefully!

Tune the probability of the agent.

$$q_{t,j} \uparrow \text{ if } \hat{r}_{t,j} \uparrow$$

Choose your **update rule** very carefully!

How to estimate and aggregate?

$$\forall t \geq 1$$

How to estimate and aggregate?

$$\forall t \geq 1$$

💡 Turn lasso into a sparse online regression oracle

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

$$\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x}))$$

$$\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M) \in \mathbb{R}^{dM}$$

How to estimate and aggregate?

$$\forall t \geq 1$$

💡 Turn lasso into a sparse online regression oracle

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

$$\begin{aligned} \phi(\mathbf{x}) &= (\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})) \\ \boldsymbol{\theta} &= (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M) \in \mathbb{R}^{dM} \end{aligned}$$

Theorem (Anytime Lasso Conf Seq)

For appropriate choice of $(\lambda_t)_{t \geq 1}$,



$$\mathbb{P} \left(\forall t \geq 1 : \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_t\|_2 \lesssim \sqrt{\frac{\log(M/\delta)}{t}} \right) \geq 1 - \delta$$

How to estimate and aggregate?

$$\forall t \geq 1$$

💡 Turn lasso into a sparse online regression oracle

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

$$\begin{aligned} \boldsymbol{\phi}(\mathbf{x}) &= (\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})) \\ \boldsymbol{\theta} &= (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M) \in \mathbb{R}^{dM} \end{aligned}$$

Theorem (Anytime Lasso Conf Seq)

For appropriate choice of $(\lambda_t)_{t \geq 1}$,



$$\mathbb{P} \left(\forall t \geq 1 : \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_t\|_2 \lesssim \sqrt{\frac{\log(M/\delta)}{t}} \right) \geq 1 - \delta$$

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

average reward of agent j

How to estimate and aggregate?

$$\forall t \geq 1$$

💡 Turn lasso into a sparse online regression oracle

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

$$\begin{aligned} \boldsymbol{\phi}(\mathbf{x}) &= (\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})) \\ \boldsymbol{\theta} &= (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M) \in \mathbb{R}^{dM} \end{aligned}$$

Theorem (Anytime Lasso Conf Seq)

For appropriate choice of $(\lambda_t)_{t \geq 1}$,



$$\mathbb{P} \left(\forall t \geq 1 : \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_t\|_2 \lesssim \sqrt{\frac{\log(M/\delta)}{t}} \right) \geq 1 - \delta$$

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

average reward of agent j

💡 Exponential Weighting

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

How to estimate and aggregate?

$$\forall t \geq 1$$

💡 Turn lasso into a sparse online regression oracle

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

$$\begin{aligned} \boldsymbol{\phi}(\mathbf{x}) &= (\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})) \\ \boldsymbol{\theta} &= (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M) \in \mathbb{R}^{dM} \end{aligned}$$

Theorem (Anytime Lasso Conf Seq)

For appropriate choice of $(\lambda_t)_{t \geq 1}$,



$$\mathbb{P} \left(\forall t \geq 1 : \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_t\|_2 \lesssim \sqrt{\frac{\log(M/\delta)}{t}} \right) \geq 1 - \delta$$

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

average reward of agent j

💡 Exponential Weighting

estimate of the reward obtained
by agent j so far

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

How to estimate and aggregate?

$\forall t \geq 1$

💡 Turn lasso into a sparse online regression oracle

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

$\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x}))$
 $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M) \in \mathbb{R}^{dM}$

Theorem (Anytime Lasso Conf Seq)

For appropriate choice of $(\lambda_t)_{t \geq 1}$,



$$\mathbb{P} \left(\forall t \geq 1 : \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_t\|_2 \lesssim \sqrt{\frac{\log(M/\delta)}{t}} \right) \geq 1 - \delta$$

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

average reward of agent j

💡 Exponential Weighting

estimate of the reward obtained
by agent j so far

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

sensitivity of updates

Putting it all together: ALExp

Anytime **Exponential** weighting algorithm with **Lasso** reward estimates

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

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

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

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

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

Update selection distribution

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

Putting it all together: ALExp

Anytime **Exponential** weighting algorithm with **Lasso** reward estimates

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

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

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

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

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

Update selection distribution

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

prescribed in the paper

Theorem (Online Model Selection)

For appropriate choices of parameters,

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

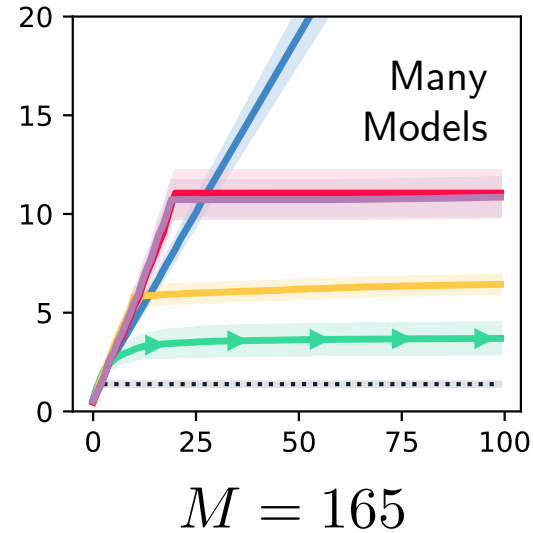
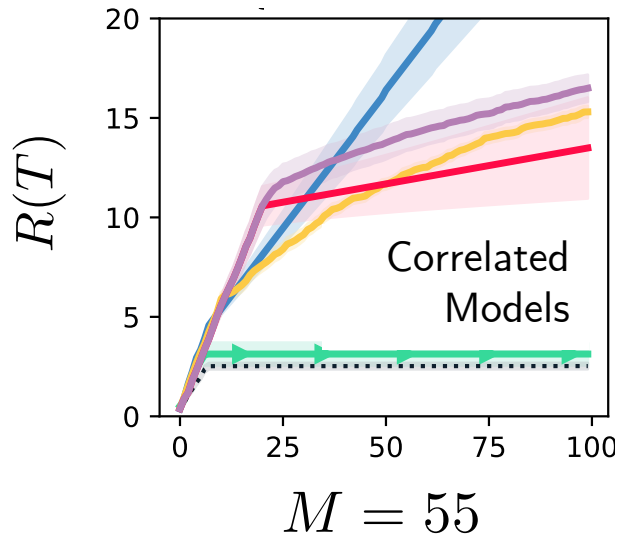
w.h.p. simultaneously for all $T \geq 1$.



Synthetic Experiments

data generation & baselines described in the paper.

..... Oracle UCB
knows j^*
uses all features
Explore then commit
[Agarwal et al. 2017]
→ ALEXP
Naive UCB
ETC
ETS
Corral



More experiments: