

---

# Anytime Model Selection in Linear Bandits

---

Parnian Kassraie<sup>1</sup> Nicolas Emmenegger<sup>1</sup> Andreas Krause<sup>1</sup> Aldo Pacchiano<sup>2,3</sup>  
<sup>1</sup>ETH Zurich <sup>2</sup>Broad Institute of MIT and Harvard <sup>3</sup>Boston University  
{pkassraie, nicolaem, krausea}@ethz.ch pacchian@bu.edu

## Abstract

Model selection in the context of bandit optimization is a challenging problem, as it requires balancing exploration and exploitation not only for action selection, but also for model selection. One natural approach is to rely on online learning algorithms that treat different models as experts. Existing methods, however, scale poorly (poly $M$ ) with the number of models  $M$  in terms of their regret. Our key insight is that, for model selection in linear bandits, we can emulate full-information feedback to the online learner with a favorable bias-variance trade-off. This allows us to develop ALEXP, which has an exponentially improved ( $\log M$ ) dependence on  $M$  for its regret. ALEXP has anytime guarantees on its regret, and neither requires knowledge of the horizon  $n$ , nor relies on an initial purely exploratory stage. Our approach utilizes a novel time-uniform analysis of the Lasso, establishing a new connection between online learning and high-dimensional statistics.

## 1 Introduction

When solving bandit problems or performing Bayesian optimization, we need to commit to a reward model *a priori*, based on which we estimate the reward function and build a policy for selecting the next action. In practice, there are many ways to model the reward by considering different feature maps or hypothesis spaces, e.g., for optimizing gene knockouts [Gonzalez et al., 2015, Pacchiano et al., 2022] or parameter estimation in biological dynamical systems [Ulmasov et al., 2016, Imani et al., 2019]. It is not known a priori which model is going to yield the most sample efficient bandit algorithm, and we can only select the right model as we gather empirical evidence. This leads us to ask, can we perform adaptive model selection, while simultaneously optimizing for reward?

In an idealized setting with no sampling limits, given a model class of size  $M$ , we could initialize  $M$  bandit algorithms (a.k.a. agents) in parallel, each using one of the available reward models. Then, as the algorithms run, at every step we can select the most promising agent, according to the cumulative rewards that are obtained so far. Model selection can then be cast into an online optimization problem on a  $M$ -dimensional probability simplex, where the probability of selecting an agent is dependent on its cumulative reward, and the optimizer seeks to find the distribution with the best return in hindsight. This approach is impractical for large model classes, since at every step, it requires drawing  $M$  different samples in parallel from the environment so that the reward for each agent is realized.

In realistic applications of bandit optimization, we can only draw *one* sample at a time, and so we need to design an algorithm which allocates more samples to the agents that are deemed more promising. Prior work [e.g., Maillard and Munos, 2011, Agarwal et al., 2017] propose to run a single meta algorithm which interacts with the environment by first selecting an agent, and then selecting an action according to the suggestion of that agent. The online model selection problem can still be emulated in this setup, however this time the optimizer receives partial feedback, coming from only one agent. Consequently, many agents need to be queried, and the overall regret scales with poly $M$ , again restricting this approach to small model classes. In fact, addressing the limited scope of such algorithms, Agarwal et al. [2017] raise an open problem on the feasibility of obtaining a  $\log M$  dependency for the regret.

We show that this rate is achievable, in the particular case of linearly parametrizable rewards. We develop a technique to “hallucinate” the reward for every agent that was not selected, and run the online optimizer with emulated full-information feedback. This allows the optimizer to assess the quality of the agents, without ever having queried them. As a result, our algorithm, ALEXP, satisfies a regret of rate  $O(\max\{n \log^3 M; n^{3/4} \sqrt{\log M}\})$ , with high probability, simultaneously for all  $T \geq 1$  (Theorem 1). Our key idea, leading to  $\log M$  dependency, is to employ the Lasso as a low-variance online regression oracle, and estimate the reward for the agents that were not chosen. This trick is made possible through our novel time-uniform analysis of online Lasso regression (Theorem 3). Consequently, ALEXP is horizon-independent, and explores adaptively without requiring an initial exploration stage. Empirically we find that ALEXP consistently outperforms prior work across a range of environments.

Table 1: Overview of literature on online model selection for bandit optimization

	MS Technique	Regret	MS Guarantee	adaptive exploration	anytime
Sparse Linear Bandits	Lasso	$\log M$	7	7	7
MS for Black-Box Bandits	OMD with bandit feedback	$\text{poly} M$	3	3	7
MS for Linear Bandits (Ours)	EXP4 with full-information	$\log M$	3	3	3

## 2 Related Work

Online Model selection (MS) for bandits considers combining a number of agents in a master algorithm, with the goal of performing as well as the best agent [Maillard and Munos, 2011, Agarwal et al., 2017, Pacchiano et al., 2020, Luo et al., 2022]. This literature operates on black-box model classes of size  $M$  and uses variants of Online Mirror Descent (OMD) to sequentially select the agents. The optimizer operates on importance-weighted estimates of the agents’ rewards, which has high variance ( $\text{poly} M$ ) and is non-zero only for the selected agent. Effectively, the optimizer receives partial (a.k.a. bandit) feedback and agents are at risk of starvation, since at every step only the selected agent gets the new data point. These works assume knowledge of the horizon, however, we suspect that this may be lifted with a finer analysis of OMD.

Sparse linear bandits use sparsity-inducing methods, often Lasso [Tibshirani, 1996], for estimating the reward in presence of many features, as an alternative to model selection. Early results are in the data-rich regime where the stopping time is known and larger than the number of models, and some suffer from  $\text{poly} M$  regret dependency [e.g., Abbasi-Yadkori et al., 2012]. Recent efforts often consider the contextual case, where at every step only a finite stochastic subset of the action domain is presented to the agent, and that the distribution of these points is i.i.d. and sufficiently diverse [Li et al., 2022, Bastani and Bayati, 2020, Kim and Paik, 2019, Oh et al., 2021, Cella and Pontil, 2021]. We do not rely on such assumptions. Most sparse bandit algorithms either start with a purely exploratory phase [Kim and Paik, 2019, Li et al., 2022, Hao et al., 2020, Jang et al., 2022], or rely on a priori scheduled exploration [Bastani and Bayati, 2020]. The exploration budget is set according to the horizon. Therefore, such algorithms inherently require the knowledge of the horizon and can be made anytime only via the doubling trick [Auer et al., 1995]. Table 2 presents an in-depth overview.

ALEXP inherits the best of both worlds (Table 1): its regret enjoys  $\log M$  dependency of sparse linear bandits even on compact domains, and it has adaptive probabilistic exploration with anytime guarantees. In contrast to prior literature, we perform model selection with an online optimizer (EXP4), which hallucinates full-information feedback using a low-variance Lasso estimator instead of importance-weighted estimates. Moreover, our anytime approach lifts the horizon dependence and the exploration requirement of sparse linear bandits.

Our work is inspired by and contributes to the field of Learning with Expert Advice, which analyzes incorporating the advice of  $M$  oblivious (non-adaptive) experts, with bandit or full-information feedback [Haussler et al., 1998, Auer et al., 2002b, McMahan and Streeter, 2009]. The idea of employing an online optimizer for learning stems from this literature, and has been used in various applications of online learning [Foster et al., 2017, Singla et al., 2018, Muthukumar et al., 2019, Karimi et al., 2021, Liu et al., 2022]. In particular, we are inspired by Foster and Rakhlin [2020] and Moradipari et al. [2022], who apply EXP4 to least squares estimates, for arm selection in a

contextual bandits. However, their algorithms are not anytime, and due to the choice of estimator, the corresponding regret scales with  $\mathcal{O}(\min(\sqrt{M}; \sqrt{K \log M}))$ .

### 3 Problem Setting

We consider a bandit problem where a learner interacts with the environment in rounds. At step  $t$ , the learner selects an action  $x_t \in X$ , where  $X \subseteq \mathbb{R}^{d_0}$  is a compact domain and observes a noisy reward  $y_t = r(x_t) + \epsilon_t$  such that  $\epsilon_t$  is an i.i.d. zero-mean sub-Gaussian variable with parameter  $\sigma$ . We assume the reward function  $r: X \rightarrow \mathbb{R}$  is linearly parametrizable by some unknown feature map, and that the model class  $\mathcal{F} = \{f_j: \mathbb{R}^{d_0} \rightarrow \mathbb{R}; j = 1, \dots, M\}$  contains the set of plausible feature maps. We consider the setting where  $M$  can be very large, and while the set  $\mathcal{F}$  may include misspecified feature maps, it contains at least one feature map that represents the reward function, i.e., there exists  $j^* \in [M]$  such that  $r(x) = f_{j^*}(x)$ . We assume the image of  $\mathcal{F}$  spans  $\mathbb{R}^d$ , and no two feature maps are linearly dependent, i.e. for any  $j_1, j_2 \in [M]$ , there exists no  $\alpha \in \mathbb{R}$  such that  $f_{j_1}(x) = \alpha f_{j_2}(x)$ . This assumption, which is satisfied by design in practice, ensures that the features are not ill-posed and we can explore in all relevant directions. We assume that the concatenated feature map  $\phi(x) := (\phi_1(x); \dots; \phi_M(x))$  is normalized  $\|\phi(x)\|_2 = 1$  for all  $x \in X$  and that  $\|\phi_j\|_2 \leq B$ , which implies  $\|\phi(x)\|_2 \leq B$  for all  $x \in X$ .

We will model this problem in the language of model selection where a meta algorithm aims to optimize the unknown reward function by relying on a number of base learners. In order to interact with the environment the meta algorithm selects an agent that in turn selects an action. In our setting we think of each of these feature maps as controlled by a base agent running its own algorithm. Base agent  $j$  uses the feature map  $\phi_j$  for modeling the reward. At step  $t$  of the bandit problem, each agent  $j$  is given access to the full history  $\lambda_{t-1}^j := (x_1; y_1; \dots; x_{t-1}; y_{t-1})$ , and uses it to locally estimate the reward as  $\hat{r}_{t-1}^j(x)$ , where  $\lambda_{t-1}^j \in \mathbb{R}^d$  is the estimated coefficients vector. The agent then uses this estimate to develop its action selection policy  $p_j(x) \in \mathcal{M}(X)$ . Here,  $\mathcal{M}$  denotes the space of probability measures defined on  $X$ . The condition on existence of  $j^*$  will ensure that there is at least one agent which is using a correct model for the reward, and thereby can solve the bandit problem if executed in isolation. We refer to agents as the oracle agent.

Our goal is to find a sample-efficient strategy for iterating over the agents, such that their suggested actions maximize the cumulative reward  $R_n$  achieved over any horizon  $n$ . This is equivalent to minimizing the cumulative regret  $R_n = \sum_{t=1}^n r(x^*) - \sum_{t=1}^n r(x_t)$ , where  $x^*$  is a global maximizer of the reward function. We neither know, nor assume knowledge of it.

### 4 Method

As warm-up, consider an example with deterministic agents, i.e., where  $p_j$  is a Dirac measure on a specific action  $x_{t,j}$ . Suppose it was practically feasible to draw the action suggested by every agent and observe the corresponding reward vector  $(r_{t,j})_{j=1}^M$ . In this case, model selection becomes a full-information online optimization problem, and we can design a minimax optimal algorithm as follows. We assign a probability distribution  $q_t = (q_{t,j})_{j=1}^M$  to the models, and update it such that the overall average return  $\sum_{t=1}^n q_t^T r_t$  is competitive to the best agent's average return  $\sum_{t=1}^n r_{t,j^*}$ . At every step, we update  $q_{t+1,j} \propto \exp(\sum_{s=1}^t r_{s,j})$ , since such exponential weighting is known to lead to an optimal solution for this classic online learning problem [Cesa-Bianchi and Lugosi, 2006]. In our setting however, the agents are stochastic, and we do not have access to the full.

We propose the Anytime EXponential weighting algorithm based on sparse reward estimates (ALEXP), summarized in Algorithm 1. At step  $t$  we first sample an agent  $j_t$ , and then sample an action  $x_t$  according to the agent's policy  $p_{j_t}$ . Let  $\mathcal{M}$  denote the  $M$ -dimensional probability simplex. We maintain a probability distribution  $q_t \in \mathcal{M}$  over the agents, and update it sequentially as we accumulate evidence on the performance of each agent. Ideally, we would have adjusted according to the average return of model  $j$  that is,  $E_{x \sim p_{j_t}} r(x)$ . However, since  $r$  is unknown, we estimate the average reward with some  $\hat{r}_{s,j}$ . We then update  $q_t$  for the next step via,

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

for all  $j = 1; \dots; M$ , where  $\eta_t$  is the learning rate, and controls the sensitivity of the updates. This rule allows us to imitate the full-information example that we mentioned above. By utilizing  $\hat{\mu}_{t,j}$  and hallucinating feedback from all agents, we can reduce the probability of selecting a badly performing agent, without ever having sampled them (c.f. Fig. 4). It remains to design the estimator  $\hat{\mu}_{t,j}$ . We concatenate all feature maps, and, knowing that many features are redundant, use a sparsity inducing estimator over the resulting coefficients vector. Mainly, let  $(\mu_1; \dots; \mu_M) \in \mathbb{R}^{Md}$  be the concatenated coefficients vector. We then solve

$$\hat{\mu}_t = \arg \min_{\mu \in \mathbb{R}^{Md}} L(\mu; H_t; r_t) = \arg \min_{\mu \in \mathbb{R}^{Md}} \frac{1}{t} \langle \mu, r_t \rangle - \frac{\lambda}{2} \sum_{j=1}^M \|\mu_j\|_2^2 \quad (1)$$

where  $H_t = [\phi_s(x_s)]_{s=1}^t \in \mathbb{R}^{t \times Md}$  is the feature matrix,  $r_t \in \mathbb{R}^t$  is the concatenated reward vector, and  $\lambda_t$  is an adaptive regularization parameter. Problem (1) is the online variant of the group Lasso [Lounici et al., 2011]. The second term is the loss is the  $\ell_2$ -norm of  $\mu$ , which can be seen as the  $\ell_1$ -norm of the vector  $(\|\mu_1\|_2; \dots; \|\mu_M\|_2) \in \mathbb{R}^M$ . This norm induces sparsity at the group level, and therefore, the sub-vectors  $\mu_j \in \mathbb{R}^{Md}$  that correspond to redundant feature maps are expected to be 0, i.e. the null vector. We then estimate the average return of each model by simply taking an expectation  $\hat{\mu}_{t,j} = E_{x \sim p_{t+1,j}}[\hat{\mu}_t(x)]$ . This quantity is the average return of the agent's policy  $p_{t+1,j}$ , according to our Lasso estimator. In Section 5.2 we explain why the particular choice of Lasso is crucial for obtaining  $\log M$  rate for the regret.

For action selection, with probability  $\eta_t$ , we sample agent  $j$  with probability  $q_{t,j}$  and draw  $x_t \sim p_{t,j}$  as per suggestion of the agent. With probability  $1 - \eta_t$ , we choose the action according to some exploratory distribution  $\nu_t \in \Delta(M)$  which aims to sample informative actions. This can be any design where  $\text{supp}(\nu_t) = X$ . We mix  $p_{t,j}$  with  $\nu_t$ , to collect sufficiently diverse data for model selection. We are not restricting the agents' policy, and therefore can not rely on them to explore adequately. In Theorem 1, we choose a decreasing sequence  $\eta_t$  of the probabilities of exploration at step  $t$ , since less exploration will be required as data accumulates. To conclude, the action selection policy of ALEXP is formally described as the mixture

$$p_t(x) = \eta_t \sum_{j=1}^M q_{t,j} p_{t,j}(x) + (1 - \eta_t) \nu_t(x)$$

## 5 Main Results

For the regret guarantee, we consider specific choices of base agents and exploratory distribution. Our analysis may be extended to include other policies, since ALEXP can be wrapped around any bandit agent that is described by some  $p_{t,j}$ , and allows for random exploration with any distribution  $\nu_t$ .

**Base Agents.** We assume that the oracle agent has either a UCB [Abbasi-Yadkori et al., 2011] or a GREEDY [Auer et al., 2002a] policy, and all other agents are free to choose any arbitrary policy. Similar treatment can be applied to cases where the oracle uses other (sublinear) policies for solving linear or contextual bandits [e.g., Thompson, 1933, Kaufmann et al., 2012, Agarwal et al., 2014]. In either case, agent  $j$  calculates a ridge estimate of the coefficients vector based on the history

$$\hat{\mu}_{t,j} := \arg \min_{\mu \in \mathbb{R}^{Md}} \langle \mu, y_{t,j} \rangle - \frac{\lambda}{2} \|\mu\|_2^2 = \sum_{s=1}^t \phi_s(x_s) \mu_j + \lambda^{-1} \sum_{s=1}^t y_{t,j} \phi_s(x_s)$$

Here  $\phi_s(x_s) \in \mathbb{R}^{Md}$  is the feature matrix, where each row is  $\phi_s(x_s)$  and  $\lambda$  is the regularization constant. Then at step  $t$  a GREEDY oracle suggests the action which maximizes the reward estimate  $\hat{\mu}_{t,j}^\top \phi_s(x_s)$ , and a UCB oracle queries  $\arg \max_{j \in [M]} \hat{\mu}_{t,j}^\top \phi_s(x_s)$  where  $u_{t,j}(\cdot)$  is the upper confidence bound that this agent calculates (for Proposition 21 shows that the sequence  $(u_{t,j})_{t=1}^T$  is in fact an anytime valid upper bound for the entire domain).

**Exploratory policy.** Performance of ALEXP depends on the quality of the samples that  $\nu_t$  suggests. The eigenvalues of the covariance matrix  $\Sigma_t := E_{x \sim \nu_t}[\phi(x) \phi(x)^\top]$  reflect how diverse the data is, and thus are a good indicator for data quality. van de Geer and Bühlmann [2009] present a survey on the notions of diversity defined based on these eigenvalues. Let  $\lambda_{\min}(A)$  denote the minimum eigenvalue of a matrix  $A$ . Similar to Hao et al. [2020], we assume that the maximizer of the problem below and present our regret bounds in terms of

$$C_{\min} = C_{\min}(X; \lambda) := \max_{\mu \in \mathbb{R}^{Md}} \min_{x \in X} \langle \mu, \phi(x) \rangle \quad (2)$$

---

**Algorithm 1** ALEXP
 

---

Inputs:  $(\epsilon; \delta; \eta)$  for  $t = 1$   
 Let  $q_1 \sim \text{Unif}(M)$  and initialize base agent  $(p_{1,1}; \dots; p_{1,M})$ .  
 for  $t = 1$  do  
   Draw  $\epsilon_t \sim \text{Bernoulli}(\epsilon)$ . . Decide to explore or exploit  
   if  $\epsilon_t = 1$  then  
     Choose action  $x_t$  randomly according to  $q_t$ . . Explore  
   else  
     Draw  $j_t \sim q_t$ . . Select an agent  
     Draw  $x_t \sim p_{t,j_t}$ . . Select the action suggested by the agent  
   end if  
   Observe  $y_t = r(x_t) + g_t$ . . Receive reward  
    $H_t = H_{t-1} \cup \{(x_t; y_t)\}$ . . Append history  
    $\hat{\mu}_{t,j_t} = \arg \min_{j \in [M]} L(\epsilon; H_t; j)$ . . Update the parameter estimate  
   Report  $H_t$  to all agents, and get updated policies  $(p_{t+1,1}; \dots; p_{t+1,M})$ . . Update agents  
   Update estimated average return of every agent  
     
$$\hat{\mu}_{t,j} = E_{x \sim p_{t+1,j}} [\hat{\mu}_t(x)]; \quad j = 1; \dots; M$$
  
   Update agent probabilities  
     
$$q_{t+1,j} = \frac{\exp(\epsilon_t \sum_{s=1}^t \hat{\mu}_{s,j})}{\sum_{i=1}^M \exp(\epsilon_t \sum_{s=1}^t \hat{\mu}_{s,i})}$$
  
 end for

---

which is greater than zero under the conditions specified in our problem setting. Prior works in the sparse bandit literature all require a similar or stronger assumption of this kind, and Table 2 gives an overview. Alternatively, one can work with an arbitrary  $g_t$ , as long as  $\min_j (g_t(x; j))$  is bounded away from zero. Appendix C.1 reviews some configurations  $(\epsilon; \delta; \eta)$  that lead to a non-zero minimum eigenvalue, and Corollary 12 bounds the regret of ALEXP with uniform exploration.

For this choice of agents and exploratory distribution, Theorem 1 presents an informal regret bound. Here, we have used the  $\tilde{O}$  notation, and only included the fastest growing terms. The inequality is made exact in Theorem 14, up to constant multiplicative factors.

**Theorem 1 (Cumulative Regret of ALEXP, Informal).** Let  $\epsilon \in (0; 1]$  and set  $\eta$  to be the maximizer of (2). Choose learning rate  $\epsilon_t = O(C_{\min} t^{-1/2} = C(M; \delta; \eta))$ , exploration probability  $\epsilon_t = O(t^{-1/4})$  and Lasso regularization parameter  $\eta = O(C_{\min} t^{-1/2} = C(M; \delta; \eta))$ , where

$$C(M; \delta; \eta) = O \left( \frac{q}{1 + \log(M)} + \frac{p}{d (\log(M) + (\log \log d)_+)} \right)$$

Then ALEXP satisfies the cumulative regret

$$\begin{aligned}
 R(n) = O & \left( n^{3/4} B + n^{3/4} C(M; \delta; \eta) + C_{\min}^{-1} \frac{p}{n} \bar{C}(M; \delta; \eta) \log M \right. \\
 & \left. + C_{\min}^{-1/2} n^{5/8} \frac{p}{d \log(n) + \log(1/\epsilon)} \right)
 \end{aligned}$$

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

In this bound, the first term is the regret incurred at exploratory steps (when  $\epsilon_t = 1$ ), the second term is due to the estimation error of Lasso (i.e.,  $\hat{\mu}_{t,j}$ ), and the third term is the regret of the exponential weights sub-algorithm. The fourth term, is the regret bound for the oracle agent run within the ALEXP framework. It does not depend on the agent's policy (greedy or optimistic), and is worse than the minimax optimal rate of  $\frac{p}{d \log n}$ . This is because the oracle is suggesting actions based on the history, which includes uninformative action-reward pairs queried by other, potentially misspecified, agents. In Corollary 12, we provide a regret bound independent of  $n$  for orthogonal feature maps, and show that the  $\frac{p}{n \log^3 M}; n^{3/4} \frac{p}{\log M}$  rate may be achieved even with the simple choice  $g_t \sim \text{Unif}(X)$ .

### 5.1 Proof Sketch

The regret is caused by two sources: selecting a sub-optimal agent, and an agent selecting a sub-optimal action. Accordingly, for any  $2 \leq j \leq M$ , we decompose the regret as

$$R(n) = \sum_{t=1}^n r(x_t^*) - r(x_t) = \sum_{t=1}^n r(x_t^*) - r_{t,j} + \sum_{t=1}^n r_{t,j} - r(x_t) : \quad (3)$$

The first term shows the cumulative regret of agent  $j$  when run within ALEXP. The second term evaluates the received reward against the cumulative average reward of  $j$ . We bound each term separately.

**Virtual Regret.** The first term  $R_j(n) := \sum_{t=1}^n r(x_t^*) - r_{t,j}$  compares the suggestion of agent  $j$  against the optimal action. We call it the virtual regret since the sequence  $\{x_t^*\}_{t=1}^n$  of the actions suggested by model  $j$  are not necessarily selected by the meta algorithm, unless  $j = 1$ . This regret is merely a technical tool, and not actually realized when running ALEXP. The virtual regret of the oracle agent may still be analyzed using standard techniques for linear bandits, e.g., [Abbasi-Yadkori et al. \[2011\]](#), however we need to adapt it to take into account a subtle difference: The confidence sequence of model  $j$  is constructed according to the true sequence of actions  $\{x_t^*\}_{t=1}^n$ , while its virtual regret is calculated based on the virtual sequence  $\{x_t^*\}_{t=1}^n$ , which the model suggests. The two sequences only match at the steps when model  $j$  is selected. Adapting the analysis of [Abbasi-Yadkori et al. \[2011\]](#) to this subtlety, we obtain in Lemma 15 that with probability greater than  $1 - \delta$ , simultaneously for all  $j = 1, \dots, M$ ,

$$R_j(n) = O\left(n^{5/8} C_{\min}^{1-2P} \sqrt{d \log(n) + \log(1/\delta)}\right) :$$

**Model Selection Regret.** The second term in (3) is the model selection regret  $R(n; j) := \sum_{t=1}^n r_{t,j} - r(x_t)$ , which evaluates the chosen action by ALEXP against the suggestion of the  $j$ th agent. Our analysis relies on a careful decomposition  $R(n; j) =$

$$R(n; j) = \sum_{t=1}^n \left[ \underbrace{r_{t,j} - \hat{r}_{t,j}}_{(I)} + \underbrace{\hat{r}_{t,j} - \hat{r}_{t,j}}_{(II)} + \underbrace{\hat{r}_{t,j} - \mathbb{E}[\hat{r}_{t,j}]}_{(III)} + \underbrace{\mathbb{E}[\hat{r}_{t,j}] - r(x_t)}_{(IV)} \right] :$$

We bound away each term in a modular manner, until we are left with the regret of the standard exponential weights algorithm. The terms (I) and (III) are controlled by the bias of the Lasso estimator, and are  $O(n^{3/4} C(M; \delta; d))$  (Lemma 19). The last term (IV) is zero in expectation, and reflects the deviation of  $r(x_t)$  from its mean. We observe that the summands form a bounded Martingale difference sequence, and their sum grows as  $\sqrt{n}$  (Lemma 18). Term (II) is the regret of our online optimizer, which depends on the variance of the Lasso estimator. We bound this term with  $O(n^{3/4} C(M; \delta; d) \log M)$ , by first conducting a standard anytime analysis of exponential weights (Lemma 17), and then incorporating the anytime Lasso variance bound (Lemma 20). We highlight that neither of the above steps require assumptions about the base agents. Combining these steps, Lemma 16 establishes the formal bound on the model selection regret.

**Anytime Lasso.** We develop novel confidence intervals for Lasso with history-dependent data, which are uniformly valid over an unbounded time horizon. This result may be of independent interest in applications of Lasso for online learning or sequential decision making. Here, we use these confidence intervals to bound the bias and variance terms that appear in our treatment of the model selection regret. The width of the Lasso confidence intervals depends on the quality of feature matrix  $X_t$ , often quantified by the restricted eigenvalue property [[Bickel et al., 2009](#), [van de Geer and Bühlmann, 2009](#), [Javanmard and Montanari, 2014](#)]:

**Definition 2.** For the feature matrix  $X_t \in \mathbb{R}^{d \times M}$  we define  $\lambda(t; s)$  for any  $1 \leq s \leq M$  as

$$\lambda(t; s) := \inf_{(J; b)} \frac{1}{t} \min_{\substack{X \\ \text{s.t. } b \in \mathbb{R}^d, \|b\|_0 \leq s}} \frac{\sum_{j \in J} \|b_j\|_2^2}{\sum_{j \in J} \|b_j\|_2^2} \quad \text{with } |J| \leq 3s$$

Our analysis is in terms of this quantity, and Lemma 8 explains the connection between  $\hat{C}_{\min}$  as defined in (2), particularly that  $(\hat{C}_{\min}; 2)$  is also positive with a high probability.

Theorem 3 (Anytime Lasso Confidence Sequences). Consider the data model  $y_t = \langle x_t, \theta \rangle + \epsilon_t$  for all  $t \geq 1$ , where  $\epsilon_t$  is i.i.d. zero-mean sub-Gaussian noise, and  $\theta$  is  $(F_t)_{t \geq 1}$ -predictable, where  $F_t := (x_1; \dots; x_t; \epsilon_1; \dots; \epsilon_t)$ . Then the solution of (1) guarantees

$$\mathbb{P}(\|\hat{\theta}_t - \theta\|_2 \leq \frac{4 \sqrt{p} \sqrt{\log t}}{\sqrt{2} (\hat{C}_{\min}; 2)} \mid \mathcal{F}_t) \geq 1 - \frac{1}{t}$$

if the regularization parameter satisfies for all  $t \geq 1$

$$\lambda_t \geq \frac{2 \sqrt{p}}{t} \left( 1 + \frac{12}{p} (\log(2M) + (\log \log d)_+) \right) + \frac{5 \sqrt{p}}{2} \sqrt{\frac{d (\log(2M) + (\log \log d)_+)}{t}}$$

Our confidence bound enjoys the same rate as Lasso with fixed (of i.i.d.) data, up to  $\log \log d$  factors. We prove this theorem by constructing a self-normalized martingale sequence based on the empirical process error  $\xi_t$ . We then apply the “stitched” time-uniform boundary of Howard et al. [2021]. Appendix B elaborates on this technique. Previous work on sparse linear bandits also include analysis of Lasso in an online setting, where  $F_t$  measurable. Cella and Pontil [2021] imitate the analysis and then apply a union bound over the time steps, which multiplies the width by  $\log n$  and requires knowledge of the horizon. Bastani and Bayati [2020] also rely on knowledge of  $n$  and employ a scalar-valued Bernstein inequality on the norm of the empirical process error, which inflates the width of the confidence sets by a factor of  $\sqrt{\log n}$ . We work directly on the  $\ell_2$ -norm, and use a curved boundary for sub-Gamma martingale sequences, which according to Howard et al. [2021] is uniformly tighter than a Bernstein bound, especially for small  $n$ .

## 5.2 Discussion

In light of Theorem 1 and Theorem 3, we discuss some properties of ALE

**Sparse EXP4.** Our approach presents a new connection between online learning and high-dimensional statistics. The rule for updating the probabilities in EXP4 is inspired by the exponential weighting for Exploration and Exploitation with Expert Advice (EXP4) algorithm, which was proposed by Auer et al. [2002b] and has been extensively studied in the adversarial bandit and learning with expert advice literature [e.g., McMahan and Streeter, 2009, Beygelzimer et al., 2009]. EXP4 classically uses importance-weighted (IW) or ordinary least squares (LS) estimators to estimate the reward, both of which are unbiased but high-variance choices [Bubeck et al., 2012]. In particular, in our linearly parametrized setting, the variance of IW and LS scales with  $M$ , which will lead to a  $\text{poly}(M)$  regret. However, it is known that introducing bias can be useful if it reduces the variance [Zimmert and Lattimore, 2022]. For instance, EXP3-IX [Kocák et al., 2014] and EXP4-IX [Neu, 2015] construct a biased IW estimator. Equivalently, others craft regularizers for the reward of the online optimizer, seeking to improve the bias-variance balance [e.g., Abernethy et al., 2008, Bartlett et al., 2008, Abernethy and Rakhlin, 2009, Bubeck et al., 2017, Lee et al., 2020, Zimmert and Lattimore, 2022]. A key technical observation in this work is that our online Lasso estimator can achieve sublinear regret which depends logarithmically on  $M$ . This is due to the fact that while the estimator itself is  $Md$ -dimensional, its bias squared and variance scale with  $\log M$ . To the best of our knowledge, this work is the first to instantiate the EXP4 algorithm with a sparse low-variance estimator.

**Adaptive and Anytime.** To estimate the reward, prior work on sparse bandits commonly emulate the Lasso analysis on of i.i.d. data or on a martingale sequence with a known length [Hao et al., 2020, Bastani and Bayati, 2020]. These works require a long enough sequence of exploratory samples, and knowledge of the horizon to plan this sequence. ALE EXP4 removes both of these constraints, and presents a fully adaptive algorithm. Crucially, we employ the elegant martingale bounds of Howard et al. [2020] to present the first time-uniform analysis of the Lasso with history-dependent data (Theorem 3). This way we can use all the data points and explore with a probability which vanishes at  $\mathcal{O}(t^{-4})$  rate. Our anytime confidence bound for Lasso, together with the horizon-independent analysis of the exponential weights algorithm, also allows EXP4 to be stopped at any time with valid guarantees.

**Rate Optimality.** For  $M \leq n$ , we obtain  $\mathcal{O}(\sqrt{p n \log^3 M})$  regret, which matches the rate conjectured by Agarwal et al. [2017]. However, if  $M$  is comparable to or smaller, our regret scales with

Figure 1: ALEXP can model-select in both orthogonal and correlated classes ( $M = 55$ ) Figure 2: ALEXP performs well on a large class ( $M = 165$ )

$O(n^{3-4^p} \overline{\log M})$ , and while it is still sublinear and scales logarithmically with the dependency on  $n$  is sub-optimal. This may be due to the conservative nature of our model selection analysis, during which we do not make assumptions about the dynamics of the base agents. Therefore, to ensure sufficiently diverse data for successful model selection, we need to occasionally choose exploratory actions with a vanishing probability of  $\epsilon$ . We conjecture that this is avoidable, if we make more assumptions about the agents, e.g., that a sufficient number of agents can achieve sublinear regret if executed in isolation. Banerjee et al. [2023] show that the data collected by sublinear algorithms organically satisfies a minimum eigenvalue lowerbound, which may also be sufficient for model selection. We leave this as an open question for future work.

## 6 Experiments

**Experiment Setup.** We create a synthetic dataset based on our data model (Section 3), and choose the domain to be  $d$ -dimensional  $X = [-1; 1]$ . As a natural choice of features, we consider the set of degree  $p$  Legendre polynomials, since they form an orthonormal basis for  $\mathcal{P}(X)$  if  $p$  grows unboundedly. We construct each feature map, by choosing different polynomials from this set, and therefore obtaining  $M = \binom{p+1}{s}$  different models. More formally, we let  $f_j(x) = (P_{j_1}(x); \dots; P_{j_s}(x)) \in \mathbb{R}^s$  where  $j_1; \dots; j_s \in \{0; \dots; p\}$  and  $P_{j_i}$  denotes a degree  $j_i$  Legendre polynomial. To construct the reward function, we randomly sample from  $[M]$ , and draw  $j_i$  from an i.i.d. standard gaussian distribution. We then normalize  $j_i$  by  $\|j\|$ . When sampling from the reward, we add Gaussian noise with standard deviation  $\sigma = 0.01$ . Figure 5 in the appendix shows how the random reward functions may look. For all experiments we set  $n_0 = 100$ , and plot the cumulative regret  $R(n)$  averaged over 20 different random seeds, the shaded areas in all figures show the standard error across these runs.

**Algorithms.** We perform experiments on two UCB algorithms, one with oracle knowledge of  $f^*$ , and a naive one which takes into account all feature maps. We run Explore-then-Commit (ETC) by Hao et al. [2020], which explores for a horizon of steps, performs Lasso once, and then selects actions greedily for the remaining steps. As another baseline, we introduce Explore-then-Select (ETS) that explores for  $n_0$  steps, performs model selection using the sparsity pattern of the Lasso estimator. Performance of ETC and ETS depends highly on  $n_0$ , so we tune this hyperparameter separately for each experiment. We also run CORRAL as proposed by Agarwal et al. [2017], with UCB agents similar to ALEXP. We tune the hyper-parameters of CORRAL as well. To initialize ALEXP we set the rates of  $\epsilon_t$ ;  $\delta_t$  and  $\eta_t$  according to Theorem 1, and perform a light hyper-parameter tuning to choose the scaling constants. We have included the details and results of our hyper-parameter tuning in Appendix F.1. To solve (1), we use CELER, a fast solver for the group Lasso [Massias et al., 2018]. Every time a UCB policy is used, we set the exploration coefficient  $\gamma = 2$ , and every time exploration is required, we sample according to  $\mu = \text{Unif}(X)$ . Appendix F includes the pseudo-code for all baseline algorithms.

**Easy vs. Hard Cases.** We construct an easy problem instance, where  $d = 2$ ,  $p = 10$ , and thus  $M = 55$ . Models are lightly correlated since each two model can have at most one Legendre polynomial in common. We also generate an instance with highly correlated feature maps where  $s = 8$  and  $p = 10$ , which will be a harder problem, since out of the total  $M = 55$  models, there are 36 models which have at least 6 Legendre polynomials in common with the oracle model. Figure 1 shows that not only ALEXP is not affected by the correlations between the models, but also it achieves

<sup>1</sup>The PYTHON code for reproducing the experiments is accessible on [github.com/lasgroup/ALEXP](https://github.com/lasgroup/ALEXP).



Figure 3: ALEXP is hardly affected by increasing the number of models (y-axis have various scales)

Figure 4: ALEXP can rule out models without ever having queried them ( $M = 165$ )

a performance competitive to the oracle in both cases, implying that our exponential weights technique for model selection is robust to choice of features. ETC and ETS rely on Lasso for model selection, which performs poorly in the case of correlated features. CSRRAL uses log-barrier-OMD with an importance-weighted estimator, which has a significantly high variance. The curves for CSRRAL in Figures 1 and 2 is cropped since the regret values get very large. Figure 6 shows the complete results. We construct another hard instance (Fig. 2), where the model class is large,  $p = 10$ ;  $M = 165$ . ALEXP continues to outperform all baselines with a significant gap. It is clear in the regret curves how explore-then-commit style algorithms are inherently horizon-dependent, and may exhibit linear regret, if stopped at an arbitrary time. This is not an issue with the other algorithms.

Scaling with  $M$ . Figure 3 shows how well the algorithms scale as  $M$  grows. For this experiment we set  $p = 2$  and change  $M$  from 9 to 139. While increasing  $M$  hardly affects ALEXP and Oracle UCB, other baselines become less and less sample efficient.

Learning Dynamics of ALEXP. Figure 4 gives some insight into the dynamics of ALEXP when  $M = 165$ . In particular, it shows how ALEXP can rule out sub-optimal agents without ever having queried them. Figure (a) shows the distribution at  $t = 20$  which is roughly equal to the optimal  $p_j$  for ETC in this configuration. The oracle model  $j^*$  is annotated with a star, and has the highest probability of selection. We observe that, already at this time step, more than 80% of the agents are practically ruled out, due to small probability of selection. However, according to Figure (b), which shows  $M_t$ , the total number of visited models, less than 10% of the models are queried at  $t = 20$ . This is the key practical benefit of ALEXP compared to black-box algorithms such as CSRRAL. Lastly, Figure (c) shows how  $q_{t,j^*}$ , the probability of selecting the oracle agent changes with time. While this probability is higher than that of the other agents, Figure (c) shows that it is not exceeding 0.25, therefore there is always a probability of greater than 0.75 that we sample another agent, making ALEXP robust to hard problem instances where many agents perform efficiently. We conclude that ALEXP seems to rapidly recognize the better performing agents, and select among them with high probability.

## 7 Conclusion

We proposed ALEXP, an algorithm for simultaneous online model selection and bandit optimization. As a first, our approach leads to anytime valid guarantees for model selection and bandit regret, and does not rely on a priori determined exploration schedule. Further, we showed how the Lasso can be used together with the exponential weights algorithm to construct a low-variance online learner. This new connection between high-dimensional statistics and online learning opens up avenues for future research on high-dimensional online learning. We established empirically that ALEXP has favorable exploration-exploitation dynamics, and outperforms existing baselines. We tackled the open problem of Agarwal et al. [2017], and showed that  $M$  dependency for regret is achievable for linearly parametrizable rewards. This problem remains open for more general, non-parametric reward classes.

## Acknowledgments

We thank Jonas Rothfuss and Miles Wang-Henderson for their valuable suggestions regarding the writing. We thank Scott Sussex and Felix Schur for their thorough feedback. This research was supported by the European Research Council (ERC) under the European Union's Horizon 2020 research and Innovation Program Grant agreement no. 815943. Nicolas Emmenegger was supported by the Zeno Karl Schindler Foundation and the Swiss Study Foundation. Aldo Pacchiano would like to thank the support of the Eric and Wendy Schmidt Center at the Broad Institute of MIT and Harvard. This work was supported in part by funding from the Eric and Wendy Schmidt Center at the Broad Institute of MIT and Harvard.

## References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in Neural Information Processing Systems*, 2011.
- Yasin Abbasi-Yadkori, David Pal, and Csaba Szepesvari. Online-to-con dence-set conversions and application to sparse stochastic bandits. *Arti cial Intelligence and Statistics* PMLR, 2012.
- Jacob Abernethy and Alexander Rakhlin. Beating the adaptive bandit with high probability. *Information Theory and Applications Workshop* IEEE, 2009.
- Jacob Abernethy, Elad Hazan, and Alexander Rakhlin. An ef cient algorithm for bandit linear optimization. In *Conference on Learning Theory* 2008.
- Alekh Agarwal, Daniel Hsu, Satyen Kale, John Langford, Lihong Li, and Robert Schapire. Taming the monster: A fast and simple algorithm for contextual bandits. *International Conference on Machine Learning* PMLR, 2014.
- Alekh Agarwal, Haipeng Luo, Behnam Neyshabur, and Robert E Schapire. Corraling a band of bandit algorithms. In *Conference on Learning Theory* PMLR, 2017.
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. Gambling in a rigged casino: The adversarial multi-armed bandit problem. *Foundations of Computer Science* IEEE, 1995.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning* 2002a.
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing* 2002b.
- Debangshu Banerjee, Avishek Ghosh, Sayak Ray Chowdhury, and Aditya Gopalan. Exploration in linear bandits with rich action sets and its implications for inference. *International Conference on Arti cial Intelligence and Statistics* PMLR, 2023.
- Peter L Bartlett, Varsha Dani, Thomas P Hayes, Sham M Kakade, Alexander Rakhlin, and Ambuj Tewari. High-probability regret bounds for bandit online linear optimization. *Conference on Learning Theory* 2008.
- Hamsa Bastani and Mohsen Bayati. Online decision making with high-dimensional covariates. *Operations Research* 2020.
- Alina Beygelzimer, John Langford, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandit algorithms with supervised learning guarantees. *International Conference on Arti cial Intelligence and Statistics* 2011.
- Peter J. Bickel, Ya'acov Ritov, and Alexandre B. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *The Annals of Statistics* 2009.
- Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning* 2012.
- Sébastien Bubeck, Yin Tat Lee, and Ronen Eldan. Kernel-based methods for bandit convex optimization. In *ACM SIGACT Symposium on Theory of Computing* pages 72–85, 2017.

- Peter Bühlmann and Sara Van De Geer. *Statistics for high-dimensional data: methods, theory and applications* Springer Science & Business Media, 2011.
- Alexandra Carpentier and Rémi Munos. Bandit theory meets compressed sensing for high dimensional stochastic linear bandit. *Artificial Intelligence and Statistics* PMLR, 2012.
- Leonardo Cella and Massimiliano Pontil. Multi-task and meta-learning with sparse linear bandits. In *Conference on Uncertainty in Artificial Intelligence* 2021.
- Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games* Cambridge university press, 2006.
- Dylan Foster and Alexander Rakhlin. Beyond ucb: Optimal and efficient contextual bandits with regression oracles. *International Conference on Machine Learning* PMLR, 2020.
- Dylan J Foster, Satyen Kale, Mehryar Mohri, and Karthik Sridharan. Parameter-free online learning via model selection *Advances in Neural Information Processing Systems* 2017.
- Dylan J Foster, Akshay Krishnamurthy, and Haipeng Luo. Model selection for contextual bandits. *Advances in Neural Information Processing Systems* 2019.
- Sébastien Gerchinovitz. Sparsity regret bounds for individual sequences in online linear regression. In *Conference on Learning Theory* 2011.
- Apostolos Giannopoulos. Notes on isotropic convex bodies, October 2003.
- Javier Gonzalez, Joseph Longworth, David C James, and Neil D Lawrence. Bayesian optimization for synthetic gene design *arXiv preprint* 2015.
- Botao Hao, Tor Lattimore, and Mengdi Wang. High-dimensional sparse linear bandits *Advances in Neural Information Processing Systems* 2020.
- David Haussler, Jyrki Kivinen, and Manfred K Warmuth. Sequential prediction of individual sequences under general loss functions *IEEE Transactions on Information Theory* 1998.
- Steven R. Howard, Aaditya Ramdas, Jon McAuliffe, and Jasjeet Sekhon. Time-uniform chernoff bounds via nonnegative supermartingales *Probability Surveys* 2020.
- Steven R Howard, Aaditya Ramdas, Jon McAuliffe, and Jasjeet Sekhon. Time-uniform, nonparametric, nonasymptotic confidence sequences *The Annals of Statistics* 2021.
- Mahdi Imani, Seyede Fatemeh Ghoreishi, Douglas Allaire, and Ulisses M Braga-Neto. Mfbo-ssm: Multi-fidelity bayesian optimization for fast inference in state-space models *Proceedings of the AAAI Conference on Artificial Intelligence* 2019.
- Kyoungseok Jang, Chicheng Zhang, and Kwang-Sung Jun. Popart: Efficient sparse regression and experimental design for optimal sparse linear bandits *Advances in Neural Information Processing Systems* 2022.
- Adel Javanmard and Andrea Montanari. Confidence intervals and hypothesis testing for high-dimensional regression *The Journal of Machine Learning Research* 2014.
- Mohammad Reza Karimi, Nezihe Merve Gürel, Bojan Karlaš, Johannes Rausch, Ce Zhang, and Andreas Krause. Online active model selection for pre-trained classifiers *International Conference on Artificial Intelligence and Statistics* PMLR, 2021.
- Parnian Kassarai, Jonas Rothfuss, and Andreas Krause. Meta-learning hypothesis spaces for sequential decision-making. *International Conference on Machine Learning* PMLR, 2022.
- Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On bayesian upper confidence bounds for bandit problems. In *Artificial intelligence and statistics* PMLR, 2012.
- Gi-Soo Kim and Myunghee Cho Paik. Doubly-robust lasso bandits *Advances in Neural Information Processing Systems* 32, 2019.

- Tomáš Kocák, Gergely Neu, Michal Valko, and Rémi Munos. Efficient learning by implicit exploration in bandit problems with side observations. *Advances in Neural Information Processing Systems*, 2014.
- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- Chung-Wei Lee, Haipeng Luo, Chen-Yu Wei, and Mengxiao Zhang. Bias no more: high-probability data-dependent regret bounds for adversarial bandits and *Advances in Neural Information Processing Systems*, 2020.
- Wenjie Li, Adarsh Barik, and Jean Honorio. A simple unified framework for high dimensional bandit problems. *International Conference on Machine Learning*, 2022.
- Xuefeng Liu, Fangfang Xia, Rick L Stevens, and Yuxin Chen. Cost-effective online contextual model selection. *arXiv preprint*, 2022.
- Karim Lounici, Massimiliano Pontil, Sara Van De Geer, and Alexandre B Tsybakov. Oracle inequalities and optimal inference under group sparsity. *The annals of statistics*, 2011.
- Haipeng Luo, Mengxiao Zhang, Peng Zhao, and Zhi-Hua Zhou. Corraling a larger band of bandits: A case study on switching regret for linear bandits. *Conference on Learning Theory*, 2022.
- Odalric-Ambrym Maillard and Rémi Munos. Adaptive bandits: Towards the best history-dependent strategy. *International Conference on Artificial Intelligence and Statistics*, 2011.
- Mathurin Massias, Alexandre Gramfort, and Joseph Salmon. Celer: a fast solver for the lasso with dual extrapolation. *International Conference on Machine Learning*, 2018.
- H. Brendan McMahan and Matthew Streeter. Tighter bounds for multi-armed bandits with expert advice. *Proceedings of the 22nd Annual Conference on Learning Theory*, 2009.
- Ahmadreza Moradipari, Berkay Turan, Yasin Abbasi-Yadkori, Mahnoosh Alizadeh, and Mohammad Ghavamzadeh. Feature and parameter selection in stochastic linear bandits. *International Conference on Machine Learning*, 2022.
- Vidya Muthukumar, Mitas Ray, Anant Sahai, and Peter Bartlett. Best of many worlds: Robust model selection for online supervised learning. *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, 2019.
- Gergely Neu. Explore no more: Improved high-probability regret bounds for non-stochastic bandits. *Advances in Neural Information Processing Systems*, 2015.
- Min-hwan Oh, Garud Iyengar, and Assaf Zeevi. Sparsity-agnostic lasso bandits. *International Conference on Machine Learning*, 2021.
- Aldo Pacchiano, My Phan, Yasin Abbasi-Yadkori, Anup Rao, Julian Zimmert, Tor Lattimore, and Csaba Szepesvári. Model selection in contextual stochastic bandit problems. *Advances in Neural Information Processing Systems*, 2020.
- Aldo Pacchiano, Drausin Wulsin, Robert A Barton, and Luis Voloch. Neural design for genetic perturbation experiments. *arXiv preprint*, 2022.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- Ali Rahimi, Benjamin Recht, et al. Random features for large-scale kernel machines. *Advances in Neural Information Processing Systems*, 2007.
- Felix Schur, Parnian Kassraie, Jonas Rothfuss, and Andreas Krause. Lifelong bandit optimization: No prior and no regret. *Conference on Uncertainty in Artificial Intelligence*, 2023.

- Adish Singla, Hamed Hassani, and Andreas Krause. Learning to interact with learning agents. In AAAI Conference on Artificial Intelligence 2018.
- William R Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika* 1933.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)* 1996.
- Doniyor Ulmasov, Caroline Baroukh, Benoit Chachuat, Marc Peter Deisenroth, and Ruth Misener. Bayesian optimization with dimension scheduling: Application to biological systems. *Computer Aided Chemical Engineering*. Elsevier, 2016.
- Sara A. van de Geer and Peter Bühlmann. On the conditions used to prove oracle results for the Lasso. *Electronic Journal of Statistics* 2009.
- Roman Vershynin. High-dimensional probability: An introduction with applications in data science volume 47. Cambridge university press, 2018.
- Martin J Wainwright. High-dimensional statistics: A non-asymptotic viewpoint. Cambridge university press, 2019.
- Julian Zimmert and Tor Lattimore. Return of the bias: Almost minimax optimal high probability bounds for adversarial linear bandits. *Conference on Learning Theory*. PMLR, 2022.

## Contents of Appendix

A	Extended Literature Review	14
B	Time Uniform Lasso Analysis	15
C	Results on Exploration	18
C.1	ALEXP with Uniform Exploration	20
C.2	Proof of Results on Exploration	20
D	Proof of Regret Bound	23
D.1	Proof of Model Selection Regret	24
D.2	Proof of Virtual Regret	30
E	Time-Uniform Concentration Inequalities	32
F	Experiment Details	34
F.1	Hyper-Parameter Tuning Results	34

## A Extended Literature Review

The sparse linear bandit literature considers linear reward functions of the form  $r(x) = \theta^T x$ , where  $x \in \mathbb{R}^p$ , however a sub-vector of size  $d$  is sufficient to span the reward function. This can be formulated as model selection among  $M = \binom{p}{d}$  different linear parametrizations, where each is a  $d$ -dimensional feature map. We present the bounds in terms of  $dM$  for coherence with the rest of the text, assuming that  $M = O(p)$ , which is the case when  $d \ll p$ .

Table 2 compares recent work on sparse linear bandits based on a number of important factors. In this table, the ETC algorithms follow the general format of exploring, performing parameter estimation once at  $t = n_0$ , and then repeatedly suggesting the same action which maximizes  $\hat{\mu}_t(x)$ . Explore-then-() Greedy takes a similar approach, however it does not settle on rather it continues to update the parameter estimate and selects  $a_t = \arg \max_{\hat{\mu}_t} (x)$ . The UCB algorithms iteratively update upper confidence bound, and choose actions which maximize them. The regret bounds in Table 2 are simplified to the terms with largest rate of growth, the reader should check the corresponding papers for rigorous results. Some of the mentioned bounds depend on problem-dependent parameters (e.g.  $\alpha_K$ ), which may not be treated as absolute constants and have complicated forms. To indicate such parameters we use in Table 2, following the notation of Hao et al. [2020]. Note that  $\alpha_K$  varies across the rows of the table, and is just an indicator for existence of other terms.

Abbasi-Yadkori et al. [2012] use the QSEW online regression oracle [Gerchinovitz, 2011] for estimating the parameter vector, together with UCB policy. The regression oracle is an exponential weights algorithm, which runs on the squared error loss. This subroutine, and thereby the algorithm proposed by Abbasi-Yadkori et al. [2012] are not computationally efficient, and this is believed to be unavoidable. This work considers the data-rich regime and shows  $\text{SR}(n) = O(\sqrt{dMn})$ , matching the lower bound of Theorem 24.3 in Lattimore and Szepesvári [2020].

Carpentier and Munos [2012] assume that the action set is a Euclidean ball, and that the noise is directly added to the parameter vector, i.e.  $r_t = x_t^T (\theta + \epsilon_t)$ . Roughly put, linear bandits with parameter noise are “easier” to solve than stochastic linear bandits with reward noise, since the noise is scaled proportionally to the features and does less “damage” [Chapter 29.3 Lattimore and Szepesvári, 2020]. In this setting, Carpentier and Munos [2012] present a  $\text{SR}(n)$  regret bound.

Recent work considers contextual linear bandits, where at every step a stochastic finite subset of size  $K$  from  $A$ , is presented to the agent. It is commonly assumed that members are i.i.d., and the sampling distribution is diverse and time-independent. The diversity assumption is often in the

form of a restricted eigenvalue condition (Definition 2) on the covariance of the context distribution [e.g. in, Kim and Paik, 2019, Bastani and Bayati, 2020]. Li et al. [2022] require a stronger condition which directly assumes that  $\lambda_{\min}(\Sigma_t)$  the minimum eigenvalue of the empirical covariance matrix is lower bounded. This is generally not true, but may hold with high probability. Hao et al. [2020] assume that the action set spans  $\mathbb{R}^{dM}$ . We believe that this assumption is the weakest in the literature, and conjecture that it is necessary for model selection. If not met, the agent can not explore in all relevant directions, and may not identify the relevant features. Our diversity assumption is similar to Hao et al. [2020], adapted to our problem setting. Mainly, we consider reward functions which are linearly parametrizable, i.e.  $\langle x, \theta \rangle$ , as oppose to linear rewards, i.e.  $x$ .

A key distinguishing factor between ALEXP and existing work on sparse linear bandit is that ALEXP is horizon-independent and does not rely on a forced exploration schedule. As shown on Table 2, majority of prior work relies either on an initial exploration stage, the length of which is determined according to [e.g., Carpentier and Munos, 2012, Kim and Paik, 2019, Li et al., 2022, Hao et al., 2020, Jang et al., 2022], or on a hand crafted schedule, which is again designed for a specific horizon [Bastani and Bayati, 2020]. Oh et al. [2021], which analyzes learned contextual bandits, does not require explicit exploration, and instead imposes restrictive assumptions on the diversity of context distribution, e.g. relaxed symmetry and balanced covariance. Regardless, the regret bounds hold in expectation, and are not time-uniform.

Table 2: Overview of recent work on high-dimensional Bandits. Parameters shows existence of other problem-dependent terms which are not constants, and varies across different rows. The regret bounds are simplified and are not rigorous.

	$ A_t $	data-poor	adap. exp.	any-time	action selection policy	MS algo	context or action assumpt.	Regret
Abbasi-Yadkori et al.	1	7	3	3	UCB	EXP4 on Sqrd error	A is compact	$d \sqrt{Mn}$ , w.h.p.
Foster et al.	K	3	3	7	UCB	EXP4 on Sqrd error	$\min(\cdot) > c$	$(Mn)^{3=4} K^{1=4} + d \sqrt{KdMn}$ w.h.p
Carpentier and Munos	1	3	7	7	UCB	Hard Thresh.	A is a ball param. noise	$d^p \bar{n}$ , w.h.p.
Bastani and Bayati	K	7	7	7	Explore then Greedy	Lasso	$\langle \cdot \rangle > c_K$	$Kd^2(\log n + \log M)^2$ w.h.p.
Kim and Paik	K	7	7	7	Explore then Greedy	Lasso	$\langle \cdot \rangle > c_K$	$d^p \bar{n} \log(Mn)$ , w.h.p.
Oh et al.	K	3	3	7	Greedy	Lasso	$\langle \cdot \rangle > c$ + other assumts.	$d^p \bar{n} \log(Mn)$ in expectation
Li et al.	K	3	7	7	ETC	Lasso	$\min(\cdot) > c$	$(n^2 d)^{1=3} \log M \bar{n}$ in expectation
Hao et al.	1	3	7	7	ETC	Lasso	A spans $\mathbb{R}^{dM}$ + is compact	$(nd C_{\min}^{-1})^{2=3} (\log M)^{1=3}$ w.h.p.
Jang et al.	1	3	7	7	ETC	Hard Thresh.	A $[1; 1]^{M^d}$ + spans $\mathbb{R}^{M^d}$	$(nd)^{2=3} (C_{\min}^{-1} \log M)^{1=3}$ w.h.p.
ALEXP (Ours)	1	3	3	3	Greedy or UCB	EXP4 on reward est.	$\text{Im}(\cdot_j)$ spans $\mathbb{R}^d$ A is compact	$d \sqrt{n \log M (n^{1=4} + C_{\min}^{-1} \log M)}$ w.h.p

## B Time Uniform Lasso Analysis

We start by showing that the sum of squared sub-gaussian variables is a sub-Gamma process (c.f. Definition 22).

Lemma 4 (Empirical Process is sub-Gamma) For  $t \geq 1$ , suppose  $\{v_i\}_{i=1}^t$  are a sequence conditionally standard sub-Gaussians adapted to the filtration  $\mathcal{F}_t = \sigma(v_1, \dots, v_t)$ . Let  $v_t \in \mathbb{R}^d$ , and  $Z_t := \frac{v_t}{\|v_t\|}$ . Define the processes  $S_t := \sum_{i=1}^t Z_i v_i$  and  $V_t := 4 \sum_{i=1}^t v_i^2$ . Then  $(S_t)_{t=0}^1$  is sub-Gamma with variance process  $(V_t)_{t=0}^1$  and scale parameter  $\alpha = 4 \max_{i=1}^t \|v_i\|$ .

Proof of Lemma 4. By definition [c.f. Definition 1, Howard et al., 2021]  $S_t$  is sub-Gamma if for each  $\lambda \in [0, 1-c]$ , there exists a supermartingale  $M_t(\lambda)_{t=0}^1$  w.r.t.  $\mathcal{F}_t$ , such that  $\mathbb{E} M_0 = 1$  and for all  $t \geq 1$ :

$$\exp \langle S_t, \lambda \rangle \leq \frac{2}{2(1-c)} V_t M_t(\lambda) \quad \text{a.s.}$$

We show the above holds in equality by proving that the left hand side itself, is a supermartingale w.r.t.  $F_t$ . We define,  $M_t(\cdot) := \exp\left(-\frac{1}{2} \sum_{s=1}^t \langle \nabla V_s, \nabla V_s \rangle\right)$ , therefore,

$$\begin{aligned} E[M_t | F_{t-1}] &= E \left[ \exp \left( -\sum_{s=1}^t \langle \nabla V_s, \nabla V_s \rangle \right) \middle| F_{t-1} \right] \\ &= E[M_{t-1} | F_{t-1}] E \left[ \exp \left( -\langle \nabla V_t, \nabla V_t \rangle \right) \middle| F_{t-1} \right] \\ &= M_{t-1} E \left[ \exp \left( -\langle \nabla V_t, \nabla V_t \rangle \right) \right] \end{aligned}$$

Note that  $\nabla V_t$  is  $F_{t-1}$ -measurable, conditionally centered and conditionally sub-exponential with parameter  $(2, 4)$  (c.f. [Vershynin \[2018, Lemma 2.7.6\]](#) and [Wainwright \[2019, Example 2.8\]](#)). Therefore, for  $c < 1$ ,

$$E \left[ \exp \left( -\langle \nabla V_t, \nabla V_t \rangle \right) \middle| F_{t-1} \right] \leq \exp \left( -\frac{c}{2} \langle \nabla V_t, \nabla V_t \rangle \right)$$

where the last inequality holds due to the fact that  $1 - c < 1$ . Therefore,

$$E[M_t | F_{t-1}] \leq M_{t-1} \exp \left( -\frac{c}{2} \langle \nabla V_t, \nabla V_t \rangle \right) = M_{t-1}$$

for  $t \in [0, 1/c)$ , concluding the proof.  $\square$

We now construct a self-normalizing martingale sequence based on the empirical process error term, and recognize that it is a sub-gamma process. We then employ our curved Bernstein bound Lemma 25 to control the boundary. This step will allow us to "ignore" the empirical process error term later in the lasso analysis.

Lemma 5 (Empirical Process is dominated by regularization)

$$A_j = \left\{ t \in [0, 1] : \left( \sum_{i=1}^t \langle \nabla V_i, \nabla V_i \rangle \right)^2 \leq t \right\}$$

Then, for any  $0 < \epsilon < 1$ , the event  $A = \bigcap_{j=1}^M A_j$  happens with probability  $1 - \epsilon$ , if for all  $t \in [0, 1]$ ,

$$\frac{2}{p} \left( 1 + \frac{5}{2} \frac{p}{d} \left( \log(2M) + (\log \log d)_+ \right) + \frac{12}{2} \left( \log(2M) + (\log \log d)_+ \right) \right) \leq \epsilon$$

Proof of Lemma 5. This proof includes a treatment of the empirical process similar to Lemma B.1 in [Kassraie et al. \[2022\]](#), but adapts it to our time-uniform setting. Since we have zero-mean sub-gaussian variables, as driven in Lemma 3.1 [\[Lounici et al., 2011\]](#), it holds that

$$A_j^c = \left\{ t \in [0, 1] : \frac{1}{t^2} \left\| \sum_{i=1}^t \nabla V_i \right\|_2^2 \geq \frac{2}{4} = \frac{1}{2} \right\} = \left\{ t \in [0, 1] : \frac{1}{2} \sum_{i=1}^t \langle \nabla V_i, \nabla V_i \rangle \geq t \right\}$$

where  $\nabla V_i$  are sub-gaussian variables with variance proxy  $\langle \nabla V_i, \nabla V_i \rangle$ .  $\lambda_i$  denotes the  $i$ -th eigenvalue of  $\sum_{i=1}^t \nabla V_i \nabla V_i^T$  with the concatenated vector  $\nabla V = (\nabla V_1; \dots; \nabla V_t)$ , and

$$\frac{1}{t^2} \left\| \sum_{i=1}^t \nabla V_i \right\|_2^2 = \frac{\text{tr} \left( \sum_{i=1}^t \nabla V_i \nabla V_i^T \right)}{t^2}$$

We can apply Lemma 25 to control the probability of event  $A_j^c$  by tuning  $\epsilon$ . Mainly, for  $A_j^c$  to happen with probability less than  $\epsilon$ , Lemma 25 states that the following must hold for all

$$\frac{5}{2} \frac{p}{\max \{ 4k_2^2, 1 \}} \frac{r}{n} \leq \epsilon \quad (4)$$

Recall that w.l.o.g. feature maps are bounded everywhere  $\|x_i\|_2 \leq 1$ , and  $\text{rank}(\sum_{i=1}^t \nabla V_i \nabla V_i^T) \leq d$  which allows for the following matrix inequalities,

$$\text{tr} \left( \sum_{i=1}^t \nabla V_i \nabla V_i^T \right) = \sum_{i=1}^t \langle \nabla V_i, \nabla V_i \rangle$$



$$\frac{\text{tr}(\hat{\Sigma}_{t,j}^2)}{t} \leq \frac{p}{d} + \frac{1}{t} \sum_{i=1}^M k_{j,i}^2 \hat{\alpha}_{t,j}^2$$

Therefore,

$$\|v_t\|_k = \frac{1}{t} \sum_{i=1}^M k_{j,i}^2 \hat{\alpha}_{t,j}^2; \max_{t \in \mathcal{T}} v_t = \max_{t \in \mathcal{T}} \frac{1}{t} \sum_{i=1}^M k_{j,i}^2 \hat{\alpha}_{t,j}^2 \leq 1:$$

For Eq. (4) to hold, it suffices that for all  $t \in \mathcal{T}$ ,

$$\frac{2}{p} \left( 1 + \frac{5}{2} \frac{p}{d} \sqrt{4d(\log(2M) + (\log \log d)_+) + \frac{12}{p} (\log(2M) + (\log \log d)_+)} \right) \leq 1:$$

Therefore, if  $t$  are chosen to satisfy the above inequality, event  $\mathcal{A}_j^c$  happens with probability less than  $\frac{1}{M}$ . Then by applying union bound,  $\bigcup_{j=1}^M \mathcal{A}_j^c$  happens with probability less than  $\frac{1}{M}$ .  $\square$

**Proof of Theorem 3** The theorem statement requires that the regularization parameter  $\lambda$  is chosen such that condition of Lemma 5 is met, and therefore event  $\mathcal{A}$  happens with probability  $1 - \frac{1}{M}$ . Throughout this proof, which adapts the analysis of Theorem 3.1. [Lounici et al. \[2011\]](#) to the time-uniform setting, we condition on  $\mathcal{A}$  happening, and later incorporate the probability.

**Step 1.** Let  $\hat{\alpha}_t$  be a minimizer of  $L(\alpha; H_t; t)$  and  $\alpha$  be the true coefficients vector, then  $L(\hat{\alpha}_t; H_t; t) \leq L(\alpha; H_t; t)$ . Writing out the loss and re-ordering the inequality we obtain,

$$\frac{1}{t} \|\hat{\alpha}_t\|_2^2 \leq \frac{2}{t} \alpha^T \hat{\alpha}_t + 2 \sum_{j=1}^M k_{j,2} \hat{\alpha}_{t,j}^2$$

which is often referred to as the Basic inequality [[Bühlmann and Van De Geer, 2011](#)]. By Cauchy-Schwarz and conditioned on event  $\mathcal{A}$ ,

$$\alpha^T \hat{\alpha}_t \leq \left( \sum_{j=1}^M \alpha_{t,j}^2 \right)^{1/2} \left( \sum_{j=1}^M \hat{\alpha}_{t,j}^2 \right)^{1/2} \leq \frac{t}{2} \sum_{j=1}^M \alpha_{t,j}^2 + \frac{1}{2} \sum_{j=1}^M \hat{\alpha}_{t,j}^2$$

then adding  $\frac{1}{2} \sum_{j=1}^M \hat{\alpha}_{t,j}^2$  to both sides, applying the triangle inequality, and recalling from Section 3 that  $\alpha_{t,j} = 0$  for  $j \notin \mathcal{J}_t$  gives

$$\frac{1}{t} \|\hat{\alpha}_t\|_2^2 + \sum_{j=1}^M \hat{\alpha}_{t,j}^2 \leq 2 \sum_{j=1}^M \alpha_{t,j}^2 + 2 \sum_{j=1}^M k_{j,2} \hat{\alpha}_{t,j}^2$$

$$4 \sum_{j \in \mathcal{J}_t} \hat{\alpha}_{t,j}^2 \leq 2 \sum_{j=1}^M \alpha_{t,j}^2 + 2 \sum_{j=1}^M k_{j,2} \hat{\alpha}_{t,j}^2$$

Since each term on the left hand side is positive, then each is also individually smaller than the right hand side, and we obtain,

$$\frac{1}{t} \|\hat{\alpha}_t\|_2^2 \leq 4 \sum_{j \in \mathcal{J}_t} \hat{\alpha}_{t,j}^2 \tag{5}$$

$$\sum_{j \in \mathcal{J}_t} \hat{\alpha}_{t,j}^2 \leq 3 \sum_{j \in \mathcal{J}_t} \hat{\alpha}_{t,j}^2 + \sum_{j \notin \mathcal{J}_t} \hat{\alpha}_{t,j}^2 \tag{6}$$

**Step 2.** Consider a sequence  $(c_1, \dots, c_k, \dots)$ , where  $c_1 \leq c_2 \leq \dots$ , then

$$c_k \leq \frac{1}{k} \sum_{i>k} c_i \leq \frac{1}{k} \sum_{i=1}^k c_i \tag{7}$$

Define  $J_1 = \{j \in \mathcal{J}_t : \hat{\alpha}_{t,j}^2 \geq c_k\}$  and  $J_2 = \{j \in \mathcal{J}_t : \hat{\alpha}_{t,j}^2 < c_k\}$  where

$$j^0 = \arg \max_{\substack{j \in \mathcal{J}_t \\ j \in \mathcal{J}_t^0}} \hat{\alpha}_{t,j}^2$$

For any  $J \subseteq [M]$  the complementing set is denoted  $J^c = [M] \setminus J$ . For simplicity let  $c_j = \sum_{i \in J} \hat{t}_{ij}^2$ , and let  $(k)$  denote the index of the  $k$ -th largest element of  $\{c_j : j \in J\}$ . By definition of  $J_2^c$  we have,

$$\sum_{j \in J_2^c} \hat{t}_{ij}^2 = \sum_{\substack{k > 1 \\ (k) \in J_1^c}} c_k^{(7)} = \sum_{\substack{k > 1 \\ (k) \in J_1^c}} \frac{(\sum_{i \in J_1^c} c_i)^2}{k^2} \\ = \sum_{i \in J_1^c} c_i^2 \sum_{\substack{k > 1 \\ (k) \in J_1^c}} \frac{1}{k^2} \leq 9 \sum_{j \in J_2} \hat{t}_{ij}^2; \quad (6)$$

which, in turn, gives

$$\sum_{j \in J_2} \hat{t}_{ij}^2 \leq \sum_{j=1}^M \frac{1}{10} \sum_{j \in J_2} \hat{t}_{ij}^2; \quad (8)$$

Step 3. On the other hand, due to (6), and by definition of  $J_2$  it also holds that

$$\sum_{j \in J_2^c} \hat{t}_{ij}^2 \leq 3 \sum_{j \in J_2} \hat{t}_{ij}^2;$$

From the theorem assumptions and Definition 2, we know that there exists  $(\hat{t}; 2)$ , therefore by Definition 2, the feature matrix  $\hat{t}$  satisfies,

$$\sum_{j \in J_2} \hat{t}_{ij}^2 \leq \frac{1}{2(\hat{t}; 2)} \|\hat{t}\|_2^2 \\ \stackrel{(5)}{=} \frac{1}{2(\hat{t}; 2)} \sum_{j \in J_2} \hat{t}_{ij}^2 \leq \frac{1}{2(\hat{t}; 2)} \sum_{j \in J_2} \hat{t}_{ij}^2;$$

From here, by applying (8) we get,

$$\sum_{j \in J_2} \hat{t}_{ij}^2 \leq \frac{4^p}{2(\hat{t}; 2)} \|\hat{t}\|_2^2.$$

If  $\hat{t}$  are chosen according to Lemma 5, and, in turn, the inequality above hold with probability greater than  $1 - \epsilon$ .  $\square$

## C Results on Exploration

In this section we present lower-bounds on the eigenvalues of the covariance matrix as it is later used in our regret analysis. In particular, we show that the feature matrix satisfies the restricted eigenvalue condition (Definition 2) required for valid Lasso confidence set (Theorem 3), and calculate a lower bound on  $(\hat{t}; 2)$ . The lower bound is later used by Lemma 19 and Lemma 20 to develop the model selection regret. We show this bound in three steps.

Equivalent to Definition 2, we write  $(\hat{t}; s) = \inf_{\|b\|_2 = 1} \sum_{j \in J} \hat{t}_{ij} b_j^2 = \inf_{\|b\|_2 = 1} \sum_{j \in J} \hat{t}_{ij} b_j^2$  where  $s := \sum_{j \in J} \hat{t}_{ij} b_j^2$  (9)

For simplicity in notation, we further define

$$\tilde{\lambda}(\hat{t}; s) := \min_{\|b\|_2 = 1} b^T \hat{t} b; \quad (10)$$

since  $\tilde{\lambda}(\hat{t}; s) = \sum_{j \in J} \hat{t}_{ij} b_j^2$ .

Step I. Consider the exploratory steps at which  $n = 1$ . Let  $\hat{t}^{(0)}$  be a sub-matrix of  $\hat{t}$  where only rows from exploratory steps are included. Note that  $\sum_{j \in J} \hat{t}_{ij}^{(0)}$  is a random matrix, where the number of rows are also random. We show that  $\tilde{\lambda}(\hat{t}^{(0)}; s)$  is lower bounded by  $\sum_{j \in J} \hat{t}_{ij}^{(0)}$ .

Lemma 6. Suppose  $\mathcal{I}_t$  has  $t^0$  rows. Then,

$$\lambda^2(\mathcal{I}_t; \mathbf{s}) \geq \frac{t^0}{t} \lambda^2(\mathcal{I}_t; \mathbf{s})$$

Proof of Lemma 6. Let  $\mathbf{X}^{(t)} = (\mathbf{x}_t)_{s \in \mathcal{I}_t} \in \mathbb{R}^{d \times d}$  for all  $t = 1, \dots, n$ . Note that  $\mathbf{X}^{(t)}$  is positive semi-definite by construction. We have,

$$\|\mathbf{X}^{(t)}\|_2 = \sqrt{\sum_{s \in \mathcal{I}_t} \|\mathbf{x}_s\|_2^2} = \sqrt{\sum_{s \in \mathcal{I}_t} \mathbf{x}_s^\top \mathbf{x}_s} = \sqrt{\sum_{s \in \mathcal{I}_t} \mathbf{x}_s^\top \mathbf{x}_s} = \sqrt{\sum_{s \in \mathcal{I}_t} \mathbf{x}_s^\top \mathbf{x}_s}$$

where the set  $\mathcal{I}_t$  contains the indices of the exploratory steps at which the action is selected according to  $\mathcal{I}_t$ . Therefore,

$$\begin{aligned} \lambda^2(\mathcal{I}_t; \mathbf{s}) &= \frac{1}{t} \min_{s \in \mathcal{I}_t} \|\mathbf{x}_s\|_2^2 \\ &= \min_{s \in \mathcal{I}_t} \frac{1}{t} \mathbf{x}_s^\top \mathbf{x}_s = \min_{s \in \mathcal{I}_t} \frac{1}{t} \mathbf{x}_s^\top \mathbf{x}_s \\ &= \min_{s \in \mathcal{I}_t} \frac{1}{t} \mathbf{x}_s^\top \mathbf{x}_s \end{aligned}$$

where the last inequality holds due to  $\mathbf{X}^{(s)}$  being PSD. Then we have,

$$\lambda^2(\mathcal{I}_t; \mathbf{s}) \geq \min_{s \in \mathcal{I}_t} \frac{1}{t} \mathbf{x}_s^\top \mathbf{x}_s = \frac{1}{t} \lambda^2(\mathcal{I}_t; \mathbf{s})$$

□

While the number of rows of  $\mathcal{I}_t$  is a random variable, we continue to condition on the event that  $\mathcal{I}_t$  has  $t^0$  rows, and investigate the distribution of its restricted eigenvalues.

Step II. The restricted eigenvalues of the exploratory submatrix are well bounded away from zero.

Lemma 7. Let  $\hat{\mathbf{x}}_t$  be the solution to (2), and  $\mathbf{x}_t \in \mathbb{R}^d$ . Suppose  $\mathcal{I}_t$  has  $t^0$  rows. Then for all  $\delta > 0$ ,

$$\mathbb{P}(\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2 \leq \delta) \geq 1 - \frac{80s}{t^0} \frac{P}{\log(2Md) + (\log \log 4t^0)_+}$$

where  $\mathbf{x}_t = (\mathbf{x}_s)_{s \in \mathcal{I}_t} := \mathbb{E}_{\mathbf{x}}(\mathbf{x}_s)_{s \in \mathcal{I}_t}$  and  $\sim$  is defined in (10).

Step III. Remains to combine the two above lemmas and incorporate a high probability bound on showing that it is close to  $\sum_{s=1}^t \mathbf{x}_s$ .

Lemma 8. There exist absolute constants  $C_1, C_2$  which satisfy,

$$\mathbb{P}(\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2 \leq C_1 \lambda^2(\mathcal{I}_t; \mathbf{s}) t^{1/4} + C_2 t^{5/8} \frac{P}{\log(Md) + (\log \log t)_+}) \geq 1$$

if  $t = O(t^{1/4})$ . Let  $\hat{\mathbf{x}}_t$  be the solution to (2), then it further holds that

$$\mathbb{P}(\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2 \leq C_1 C_{\min} t^{1/4} + C_2 t^{5/8} \frac{P}{\log(Md) + (\log \log t)_+}) \geq 1$$

The regret analysis of Hao et al. [2020] also relies on connecting  $\lambda^2(\mathcal{I}_t; \mathbf{s})$  to  $C_{\min}$ , and for this, they use Theorem 2.4 of Javanmard and Montanari [2014]. This theorem states that there exists a problem-dependent constant  $C_{\min}$  for which  $\lambda^2(\mathcal{I}_t; \mathbf{s}) \geq C_{\min}$  with high probability, if  $n_0$  and roughly  $n_0 = O(n^2 \log M)$ . We highlight that Lemma 8, presents a lower-bound which holds for all  $t \geq 1$ , however this comes at the cost of getting a looser lower bound than the result of Javanmard and Montanari [2014] for the larger time steps. In fact, due to the sub-optimal dependency of Lemma 8 on  $t$ , we later obtain sub-optimal dependency on the horizon for the case where  $M$  is unclear to us if this rate can be improved without assuming knowledge of  $n_0$ .

For the last lemma in this section we show that the empirical sub-matrices are also bounded away from zero. This will be required later to prove Lemma 15.

Lemma 9 (Base Model  $\min$  Bound) Assume  $\mu$  is the maximizer of Eq(2). Then, with probability greater than  $1 - \delta$ , simultaneously for all  $j = 1; \dots; M$  and  $t \geq 1$ ,

$$\min_{t_j} (\sum_{t_j} t_j) \leq C_1 C_{\min} t^{3=4} + C_2 t^{3=8} \sqrt{\log(Md) + (\log \log t)_+}$$

if  $t = O(t^{1=4})$ .

### C.1 ALEXP with Uniform Exploration

We presented our main regret bound (Theorem 1) in terms of  $C_{\min}$ , which only depends on properties of the feature maps and the action domain. We give a lower-bound on  $C_{\min}$  for a toy scenario which corresponds to the problem of linear feature selection over convex action sets.

Proposition 10 (Invertible Features) Suppose  $(x) := Ax : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is an invertible linear map, and  $X \subset \mathbb{R}^d$  is a convex body. Then,

$$C_{\min} = \frac{\min(A)}{2 \max(T)} > 0$$

where  $T$  is the transformation which maps  $X$  to an isotropic body.

The lower-bound of Proposition 10 is achieved by simply exploring via  $\text{Unif}(X)$ . Inspired by Schur et al. [2023, Lemma E.13], we show that even for non-convex action domains and orthogonal feature maps, the uniform exploration yields a constant lower-bound on restricted eigenvalues.

Proposition 11 (Orthonormal Features) Suppose  $\{j : X \rightarrow \mathbb{R}\}$  are chosen from an orthogonal basis of  $L^2(X)$ , and satisfy  $\|k_j\|_{L^2(X)} = \text{Vol}(X)^{-1/2}$ . Then there exist absolute constants  $C_1$  and  $C_2$  for which the exploration distribution  $\mu = \text{Unif}(X)$  satisfies

$$\mathbb{P}(\delta t \leq \lambda_1 \leq 2(\delta t; 2)) \geq C_1 t^{-1=4} - C_2 t^{5=8} \sqrt{\log(Md) + (\log \log t)_+} \quad \forall t \geq 1$$

The  $\delta = 1$  condition is met without loss of generality, by splitting the higher dimensional feature maps and introducing more base features, which will increase  $C_{\min}$ . Moreover, the orthonormality condition is met by orthogonalizing and re-scaling the feature maps. Basis functions such as Legendre polynomials and Fourier features [Rahimi et al., 2007] satisfy these conditions.

By invoking Proposition 11, instead of Lemma 8 in the proof of Theorem 1, we obtain the regret of ALEXP with uniform exploration.

Corollary 12 (ALEXP with Uniform Exploration) Let  $\delta \in (0; 1]$ . Suppose  $\{j : X \rightarrow \mathbb{R}\}$  are chosen from an orthogonal basis of  $L^2(X)$ , and satisfy  $\|k_j\|_{L^2(X)} = \text{Vol}(X)^{-1/2}$ . Assume the oracle agent employs a UCB or a Greedy policy as laid out in Section 5. Choose  $\bar{t} = O(1 + \sqrt{tC(M; \delta; d)})$  and  $t = O(\bar{t}^{1=4})$  and  $\bar{t} = O(C(M; \delta; d) \sqrt{\bar{t}})$ , then ALEXP with uniform exploration  $\mu = \text{Unif}(X)$  attains the regret

$$\begin{aligned} R(n) = & O(Bn^{3=4} + \sqrt{C(M; \delta; d)} \log M + B^2 \sqrt{\bar{n}} + B \sqrt{n((\log \log n B^2)_+ + \log(1 - \delta))}) \\ & + (n^{3=4} + \log n)C(M; \delta; d) + n^{5=8} \sqrt{d \log n + \log(1 - \delta)} + B^2 \end{aligned}$$

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

$$C(M; \delta; d) = O\left(\frac{1}{1 - \delta} \sqrt{\frac{1}{d(\log(M) + (\log \log d)_+) + (\log(M) + (\log \log d)_+)}}\right)$$

### C.2 Proof of Results on Exploration

As an intermediate step, we consider the restricted eigenvalue property of the empirical covariance matrix. Given  $t^0$  samples, the empirical estimate of  $\Sigma$

$$\hat{\Sigma}_{t^0} := \frac{1}{t^0} \sum_{s=1}^{t^0} (x_s - \bar{x})(x_s - \bar{x})^T \quad (11)$$

where  $x_s$  are sampled according to  $\mu$ . We show that every entry of  $\hat{\Sigma}_{t^0}$  is close to the corresponding entry in  $\Sigma$ , and later use it in the proofs of eigenvalue lemmas.

Lemma 13 (Anytime Bound for The Entries of Empirical Covariance Matrix) Let  $\hat{\Sigma}_t$  be the empirical covariance matrix corresponding to  $(\Sigma; \mathcal{X}_t)$  given  $t^0$  samples. Then,

$$P \left( d_1(\Sigma; \hat{\Sigma}_t) \leq \frac{5\sqrt{p}}{t^0} \sqrt{((\log \log 4t^0)_+ + \log(2Md))} \right)$$

where  $d_1(A; B) := \max_{i,j} |A_{ij} - B_{ij}|$ .

Proof of Lemma 13. We show the element-wise convergence of  $\hat{\Sigma}_t$  for the  $(i; j)$  entry where  $i, j = 1, \dots, dM$ . Consider the random sequence  $X_s := \frac{1}{\sqrt{p}} \sum_{i=1}^p (x_s)_i (x_s)_j$ . We show that  $X_1, \dots, X_n$  satisfies conditions of Lemma 26. We first observe that

$$E[X_s] = E[X_{1:s}] = E[X_s] = \frac{1}{p} \sum_{i=1}^p (x)_i (x)_j = 0$$

since by definition  $\frac{1}{\sqrt{p}} \sum_{i=1}^p (x)_i (x)_j > \frac{1}{\sqrt{p}} \sum_{i=1}^p (x)_i (x)_j$ . Moreover, we have normalized features  $(x)_i$   $k_1$ , therefore, each entry  $(x)_i (x)_j$  is also bounded, yielding  $\|X_s\|_2 \leq k_1$ . Then Lemma 26 implies that for all  $\epsilon > 0$ ,

$$P \left( \left| \frac{1}{t^0} \sum_{s=1}^{t^0} X_s - \frac{1}{p} \sum_{i=1}^p (x)_i (x)_j \right| \geq \frac{\epsilon}{2} \right) \leq \frac{5\sqrt{p}}{t^0} \sqrt{((\log \log 4t^0)_+ + \log(2Md))} \frac{1}{\epsilon}$$

Setting  $\epsilon = \frac{1}{2} \sqrt{p} \sqrt{((\log \log 4t^0)_+ + \log(2Md))}$  and taking a union bound over all indices concludes the proof.  $\square$

We are now ready to present the proofs to the lemmas in Appendix C.

Proof of Lemma 7. By (11) we have  $\hat{\Sigma}_t = \frac{1}{t^0} \sum_{s=1}^{t^0} X_s$ , and thereby

$$\| \hat{\Sigma}_t - \Sigma \|_F = \left\| \frac{1}{t^0} \sum_{s=1}^{t^0} X_s - \Sigma \right\|_F = \frac{1}{t^0} \sum_{s=1}^{t^0} \| X_s - \Sigma \|_F$$

Inspired by Lemma 10.1 in van de Geer and Bühlmann [2009], we show that element-wise closeness of matrices  $\hat{\Sigma}_t$  and  $\Sigma$  (c.f. Lemma 13) implies closeness in

$$\begin{aligned} \| \hat{\Sigma}_t - \Sigma \|_F &\leq \frac{1}{t^0} \sum_{s=1}^{t^0} \| X_s - \Sigma \|_F \\ &\leq \frac{1}{t^0} \sum_{s=1}^{t^0} \sqrt{\sum_{i,j} (X_s)_{ij}^2 - \sum_{i,j} \Sigma_{ij}^2} \\ &\leq \frac{1}{t^0} \sum_{s=1}^{t^0} \sqrt{\sum_{i,j} (X_s)_{ij}^2} \sqrt{p} \end{aligned}$$

where the last line holds due to Hölder's. Moreover, since  $\|X_s\|_2 \leq k_1$ , for any  $J \subseteq [dM]$  where  $|J| \leq s$  it additionally holds that  $\|X_s\|_2 \leq k_1$  and

$$\|X_s\|_2 \leq (1 + 3) k_1 \sqrt{s} \leq 4 \sqrt{s} k_1 \leq 4 \sqrt{s} k_1$$

which gives,

$$\| \hat{\Sigma}_t - \Sigma \|_F \leq \frac{1}{t^0} \sum_{s=1}^{t^0} 4 \sqrt{s} k_1 \leq \frac{4 k_1}{t^0} \sum_{s=1}^{t^0} \sqrt{s}$$

Therefore by Lemma 13,

$$\| \hat{\Sigma}_t - \Sigma \|_F \leq \frac{80 \sqrt{p}}{t^0} \sqrt{((\log \log 4t^0)_+ + \log(2Md))} \quad (12)$$

with probability greater than  $1 - \epsilon$ , simultaneously for all  $t^0 \geq 1$ .  $\square$

Proof of Lemma 8. In Lemma 6 we showed that

$$\| \hat{\Sigma}_t - \Sigma \|_F \leq \frac{t^0}{t} \| \hat{\Sigma}_{t^0} - \Sigma \|_F$$

where  $t^0$  indicates the number of rows in the exploratory sub-matrix of  $\Sigma$ . Recall that  $t^0 = \sum_{s=1}^t \mathbb{1}_{\{s \leq t^0\}}$  where  $\mathbb{1}_{\{s \leq t^0\}}$  are i.i.d Bernoulli random variables with success probability  $\frac{t^0}{t}$ . Due to Lemma 24,

$$P \left( \sum_{s=1}^t \mathbb{1}_{\{s \leq t^0\}} \geq t \frac{t^0}{t} - 2 \right) \leq 2 \quad (13)$$

where

$$t := \frac{5}{2} \frac{r}{t} \frac{(\log \log t)_+ + \log(8r)}{t}; \quad t := \sum_{s=1}^X$$

Due to Lemma 7, with probability greater than  $1 - \epsilon$  the following holds for all  $t \geq 1$

$$\begin{aligned} \mathbb{P}^2(t; 2) &\leq \frac{t^0}{t} \sim (t; 2) \frac{160^p t^{0p}}{t} \frac{r}{t} \frac{(\log \log 4t^0)_+ + \log(4Md)}{t} \\ &\leq \frac{t}{t} \sim (t; 2) \frac{160}{t} \frac{t^q}{t^2} \frac{r}{t} \frac{(\log \log (4t + t))_+ + \log(4Md)}{t} \end{aligned}$$

where the second inequality holds with probability  $1 - \epsilon$ , by incorporating (13) and taking a union bound. For the rest of the proof and to keep the calculations simple, we ignore the values of the absolute constants. We use the notation  $o(f(t))$  to show that  $f(t)$  grows much faster than  $g(t)$ . More formally, if for every constant  $c$  there exists  $t_0$ , where  $g(t) \leq cf(t)$  for all  $t \geq t_0$ . Since  $s = O(s^{1-4})$  there exists  $C$  such that  $t = Ct^{3-4}$ , then it is straightforward to observe that there exists absolute constants  $C_1, C_2$  which satisfy,

$$\begin{aligned} \mathbb{P}^2(t; 2) &\leq C_1 t^{1-4} \sim (t; 2) \frac{5t^{3-2} \sim (t; 2)^p}{2} \frac{r}{t} \frac{(\log \log t)_+ + \log(8r)}{t} \\ &\leq C_2 t^{5-8p} \frac{r}{\log(Md) + (\log \log t)_+} \leq t^{5-8p} \frac{r}{\log(Md) + (\log \log t)_+} \\ &\leq C_1 t^{1-4} \sim (t; 2) \leq C_3 t^{5-8p} \frac{r}{\log(Md) + (\log \log t)_+} \end{aligned}$$

The last inequality holds since  $t^{3-2p} \log \log t = o(t^{5-8p} \log \log t)$ . The above chain of inequalities imply that there exist absolute constants  $C_1, C_2$ , for which

$$\mathbb{P} \geq 1 - \epsilon : \mathbb{P}^2(t; 2) \leq C_1 \sim (t; 2) t^{1-4} \leq C_2 t^{5-8p} \frac{r}{\log(Md) + (\log \log t)_+} \leq 1 - \epsilon$$

If  $\epsilon$  is chosen according to (2), then  $\mathbb{P} \geq 1 - \epsilon$  yielding the lemma's second argument.  $\square$

**Proof of Lemma 9.** Fix  $j \in \{1, \dots, M\}$ , and construct the set

$$B_j = \{b \in \mathbb{R}^d : \|b\|_2 \leq 1, b_j \geq 0\} = \{b_1, \dots, b_M\}; \text{ s.t. } b_j \in \mathbb{R}^d; \|b_j\|_2 \leq 1 \text{ and } b_j^0 = 0$$

Note that  $|B_j| \leq s$ . Therefore,

$$\inf_{b \in B_j} \|b\|_2 \leq \inf_{b \in B_j} \|b\|_2 = \frac{1}{t} \mathbb{P}^-(t; s)$$

Moreover, by construction of  $B_j$  we have for all  $b \in B_j$  that  $\|b\|_2 = \|b_j\|_2$ , therefore,

$$\inf_{b \in B_j} \|b\|_2^2 = \inf_{\substack{b_j \in \mathbb{R}^d \\ \|b_j\|_2 \leq 1}} \|b_j\|_2^2 = \min_j (\|b_j\|_2^2)$$

From the above equations we conclude that  $\min_j (\|b_j\|_2^2) \leq \frac{1}{t^2} \mathbb{P}^-(t; s)$ , for all  $j = 1, \dots, M$ . Therefore, using Lemma 8 we obtain that there exists  $C_2$  such that

$$\mathbb{P} \geq 1 - \epsilon; j = 1, \dots, M : \min_j (\|b_j\|_2^2) \leq C_1 C_{\min} t^{3-4} \leq C_2 t^{3-8p} \frac{r}{\log(Md) + (\log \log t)_+} \leq 1 - \epsilon$$

$\square$

**Proof of Proposition 10.** Since  $X$  is a convex body, then there exists an invertible map  $T$  such that  $T(X)$  is an isotropic body [e.g. Proposition 1.1.1., Giannopoulos, 2003]. Then by definition,  $X \sim \text{Unif}(T(X))$  is an isotropic distribution and  $\text{Cov}(X) = I_d$  [e.g., c.f. Chapter 3.3.5 Vershynin, 2018]. Since  $T$  is linear and invertible, it may be written as  $T(x) = Ax$ , where  $A$  is an invertible matrix. Therefore,

$$\text{Cov}(T(X)) = \text{Cov}(AX) = A^T \text{Cov}(X)A = A^T \text{Cov}(X)A = A^T (I_d) A = A^T A$$

As for the minimum eigenvalue, suppose  $\lambda \in \mathbb{R}^d$  and  $\|v\|_2 = 1$ , then

$$C_{\min} = \min_j (\lambda_j) = \lambda^T A^T (I_d) A v = \lambda^T A^T A v = \lambda^T (A^T A) v = \lambda^T (T^T T) v = \frac{\|Av\|_2^2}{\|v\|_2^2} = \frac{\lambda^T (T^T T) v}{\lambda^T v} = \frac{\min(A)}{\max(T)}$$

$\square$



Lemma 16 (Any-Time Model-Selection Regret, Formal). Let  $\beta \in (0, 1]$  and  $\hat{p}$  be the maximizer of (2). Suppose  $t = O(C_{\min}^{-1} \bar{C}(M; \beta; d))$  and  $t = O(t^{1-\beta})$  and  $t = O(C(M; \beta; d) \bar{C}^{-1})$ , then there exists absolute constants  $C_1, \dots, C_5$  for which ALEXP attains the model selection regret

$$\begin{aligned} R(n; j) &\leq C_1 B n^{3-4} + C_2 \bar{C}_{\min}^{-1} C(M; \beta; d) \log M \\ &\quad + C_3 B^2 C_{\min}^{-\beta} \bar{n} + C_4 B^{\frac{1}{\beta}} \frac{1}{n ((\log \log n B^2)_+ + \log(1 - \beta))} \\ &\quad + C_5 B n^{1-4} + (n^{3-4} + \frac{\log n}{C_{\min}}) C(M; \beta; d) \left( 1 + C_{\min}^{-1} n^{3-8\beta} \frac{1}{\log(Md) + (\log \log n)_+} \right) \end{aligned}$$

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

$$C(M; \beta; d) = C_6 \frac{1}{1 + \beta} \frac{1}{d((\log(Md) + (\log \log d)_+) + (\log(Md) + (\log \log d)_+))}$$

### D.1 Proof of Model Selection Regret

Our technique for bounding the model selection regret relies on a classic horizon-independent analysis of the exponential weights algorithm, presented in Lemma 17.

Lemma 17 (Anytime Exponential Weights Guarantee). Assume  $\hat{p}_{t,j} \geq 1$  for all  $1 \leq j \leq M$  and  $t \geq 1$ . If the sequence  $(\hat{p}_{t,j})_{t \geq 1}$  is non-increasing, then for all  $t \geq 1$ ,

$$\sum_{k=1}^n \hat{r}_{t,k} - \sum_{j=1}^M \sum_{t=1}^n q_{t,j} \hat{r}_{t,j} \leq \frac{\log M}{n} + \sum_{t=1}^n \sum_{j=1}^M q_{t,j} \hat{r}_{t,j}^2$$

for any arm  $k \in [M]$ .

Proof of Lemma 17. Define  $\hat{R}_{t,j} := \sum_{s=1}^t \hat{r}_{s,j}$  to be the expected cumulative reward of agent  $j$  after  $t$  steps. We rewrite for a fixed  $t$

$$\sum_{k=1}^n \hat{r}_{t,k} - \sum_{j=1}^M \sum_{t=1}^n q_{t,j} \hat{r}_{t,j} = \sum_{k=1}^n \hat{r}_{t,k} - \sum_{j=1}^M E_j [q_t \hat{R}_{t,j}] \quad (14)$$

We focus on a single term in the second sum. For any  $j$  we have

$$\begin{aligned} E_j [q_t \hat{R}_{t,j}] &= \log(\exp(-E_j [q_t \hat{R}_{t,j}])) = \log(\exp(-E_j [q_t \hat{R}_{t,j}]^{1-t})) \\ &= \frac{1}{t} \log(\exp(-E_j [q_t \hat{R}_{t,j}])) \\ &= \frac{1}{t} \log(E_j [q_t \exp(-E_j [q_t \hat{R}_{t,j}])]) \quad (15) \end{aligned}$$

The inner expectation is over  $\hat{r}$  while the outer one is over  $\hat{r}$  and therefore has no effect. Moreover,

$$\begin{aligned} \frac{1}{t} \log E_j [q_t \exp(-E_j [q_t \hat{R}_{t,j}] + q_t \hat{R}_{t,j})] &= \frac{1}{t} \log(\exp(-E_j [q_t \hat{R}_{t,j}]) E_j [q_t \exp(q_t \hat{R}_{t,j})]) \\ &= \frac{1}{t} \log E_j [q_t \exp(-E_j [q_t \hat{R}_{t,j}])] \\ &\quad + \frac{1}{t} \log E_j [q_t \exp(q_t \hat{R}_{t,j})] \quad (16) \end{aligned}$$

where again, the expectation can be reintroduced to get the last line. Combining (15) and (16),

$$E_j [q_t \hat{R}_{t,j}] = \frac{1}{t} \log E_j [q_t \exp(-E_j [q_t \hat{R}_{t,j}] + q_t \hat{R}_{t,j})] - \frac{1}{t} \log E_j [q_t \exp(q_t \hat{R}_{t,j})] \quad (17)$$

This transformation is at the core of many exponential weight proofs [Bubeck et al., 2012, Lattimore and Szepesvári, 2020]. We first bound the first term in (17):

$$\log E_j [q_t \exp(-E_j [q_t \hat{R}_{t,j}] + q_t \hat{R}_{t,j})] = \log E_j [q_t \exp(q_t \hat{R}_{t,j})] - E_j [q_t \hat{R}_{t,j}]$$



$$\begin{aligned}
& \stackrel{(I)}{=} E_{i \sim q_t} \exp(\hat{r}_{t,i}) - 1 - E_{j \sim q_t} \hat{r}_{t,j} \\
& = E_{i \sim q_t} [\exp(\hat{r}_{t,i}) - 1 - \hat{r}_{t,i}] \\
& \stackrel{(II)}{=} E_{i \sim q_t} \frac{2}{t} \hat{r}_{t,i}^2
\end{aligned} \tag{18}$$

where in (I) we use the fact that  $\log(z) = z - 1$  and in (II) we use the fact that for  $x \in [0, 1]$ , we have  $\exp(x) = 1 + x + x^2$ , and hence  $\exp(x) - 1 - x = x^2$ . For the second term in (17), we will mirror the potential argument in [Bubeck et al. \[2012\]](#), but with a slightly different potential function. We expand the definition of  $J_t$ :

$$\begin{aligned}
\frac{1}{t} \log E_{i \sim q_t} \exp(\hat{r}_{t,i}) &= \frac{1}{t} \log \frac{\prod_{i=1}^M \exp(\hat{R}_{t,i})}{\prod_{i=1}^M \exp(\hat{R}_{t-1,i})} \\
&= \frac{1}{t} \log \frac{1}{M} \sum_{i=1}^M \exp(\hat{R}_{t,i}) + \frac{1}{t} \log \frac{1}{M} \sum_{i=1}^M \exp(\hat{R}_{t-1,i}) \\
&= J_t(\hat{r}) - J_{t-1}(\hat{r});
\end{aligned} \tag{19}$$

where we define  $J_t(\hat{r}) = \frac{1}{t} \log \frac{1}{M} \sum_{i=1}^M \exp(\hat{R}_{t,i})$ . We also define  $F_t(\hat{r}) = \frac{1}{t} \log \frac{1}{M} \sum_{i=1}^M \exp(\hat{R}_{t,i})$ . We observe the relation  $J_t(\hat{r}) = F_t(\hat{r})$ . From this, it follows that for any  $\hat{r}$ , we have  $J_t(\hat{r}) - J_{t-1}(\hat{r}) = F_t(\hat{r}) - F_{t-1}(\hat{r}) \geq 0$ , by the argument in [Bubeck et al. \[2012, Theorem 3.1\]](#) that shows  $F_t(\hat{r}) \geq F_{t-1}(\hat{r})$  for any  $\hat{r}$ .

Putting together the pieces, Now, we can bound (17) by inputting (18) and (19):

$$E_{j \sim q_t} [\hat{r}_{t,j}] \leq E_{i \sim q_t} \frac{2}{t} \hat{r}_{t,i}^2 + J_t(\hat{r}) - J_{t-1}(\hat{r})$$

With this, we rewrite (14) as

$$\sum_{t=1}^n \hat{r}_{t,k} \leq \sum_{t=1}^n E_{j \sim q_t} [\hat{r}_{t,j}] \leq \sum_{t=1}^n \hat{r}_{t,k} + \sum_{t=1}^n E_{i \sim q_t} \frac{2}{t} \hat{r}_{t,i}^2 + \sum_{t=1}^n J_t(\hat{r}) - J_{t-1}(\hat{r}) \tag{20}$$

Potential manipulation We can do an Abel transformation on the sum of potentials (20), namely obtaining

$$\sum_{t=1}^n J_t(\hat{r}) - J_{t-1}(\hat{r}) = \sum_{t=1}^{n-1} (J_t(\hat{r}) - J_{t+1}(\hat{r})) + J_n(\hat{r});$$

where we used that  $J_0(\hat{r}) = 0$ . We know  $J_t(\hat{r}) \geq 0$  and  $J_t$  is decreasing and since  $J_{t+1} \leq J_t$ , we have  $J_{t+1} - J_t \leq 0$  or  $(J_t(\hat{r}) - J_{t+1}(\hat{r})) \geq 0$ , so that for any  $x \in \mathbb{R}$

$$\begin{aligned}
\sum_{t=1}^n J_t(\hat{r}) - J_{t-1}(\hat{r}) &= \sum_{t=1}^{n-1} (J_t(\hat{r}) - J_{t+1}(\hat{r})) + J_n(\hat{r}) \\
&\stackrel{(*)}{=} \frac{\log(M)}{n} + \frac{1}{n} \log \frac{1}{M} \sum_{i=1}^M \exp(\hat{R}_{n,i}) \\
&= \frac{\log(M)}{n} + \frac{1}{n} \log \frac{1}{M} \sum_{i=1}^M \exp(\hat{R}_{n,k}) \\
&= \frac{\log(M)}{n} + \sum_{t=1}^n \hat{r}_{t,k}
\end{aligned} \tag{21}$$

where (21) follows because  $\exp$  is positive and  $\log$  is decreasing (notice that we drop  $n-1$  terms from the sum). Plugging (21) into (20), we obtain

$$\begin{aligned}
\sum_{t=1}^n \hat{r}_{t,k} &\leq \sum_{t=1}^n E_{j \sim q_t} [\hat{r}_{t,j}] \leq \sum_{t=1}^n \hat{r}_{t,k} + \sum_{t=1}^n E_{i \sim q_t} \frac{2}{t} \hat{r}_{t,i}^2 + \sum_{t=1}^n J_t(\hat{r}) - J_{t-1}(\hat{r}) \\
&\leq \sum_{t=1}^n \hat{r}_{t,k} + \sum_{t=1}^n E_{i \sim q_t} \frac{2}{t} \hat{r}_{t,i}^2 + \frac{\log(M)}{n} + \sum_{t=1}^n \hat{r}_{t,k} \\
&\leq \sum_{t=1}^n E_{i \sim q_t} \frac{2}{t} \hat{r}_{t,i}^2 + \frac{\log(M)}{n}
\end{aligned}$$

$$= \sum_{t=1}^{X^n} \sum_{j=1}^{X^M} q_{t,j} r_{t,j}^2 + \frac{\log(M)}{n}$$

□

We expressed in Section 5.2, that the model selection regret of ALEXP, is closely tied to the bias and variance of the reward estimates. The following lemma formalizes this claim.

Lemma 18. (Anytime Generic regret bound) If is picked such that  $\sum_{t=1}^n \sum_{j=1}^M q_{t,j} \geq 1$  for all  $1 \leq j \leq M$  and  $1 \leq t \leq n$  almost surely, then Algorithm 1 satisfies with probability greater than  $1 - \epsilon$ , that simultaneously for all  $1 \leq i \leq M$

$$R(n; i) \leq 2B \sum_{t=1}^{X^n} \sum_{j=1}^{X^M} q_{t,j} r_{t,j}^2 + \frac{\log M}{n} \sum_{t=1}^{X^n} \sum_{j=1}^{X^M} q_{t,j} + 10B \sqrt{\frac{P}{n} ((\log \log n B^2)_+ + \log(12))}$$

where  $r_{t,j} = \sum_{i=1}^M q_{t,i} r_{t,i}$ .

Proof of Lemma 18 Let  $\tau_t$  denote the Bernoulli random variable that is equal to 1 if at step  $t$  we select action  $i$  and 0 otherwise. At each step with  $\tau_t = 1$  ALEXP accumulates a regret of at most  $2B$ , since  $k_1 \leq B$  and  $k_2 \leq 1$ . We can decompose the regret as,

$$R(n; i) = \sum_{t=1}^{X^n} 2B \tau_t + \sum_{t=1}^{X^n} (r_{t,i} - r_t)(1 - \tau_t)$$

For the first term, by Lemma 24, we have

$$\sum_{t=1}^{X^n} 2B \tau_t \leq 2B \sum_{t=1}^{X^n} \tau_t + \frac{5P}{2} \sqrt{\frac{1}{n} ((\log \log n)_+ + \log(4))}$$

simultaneously for all  $1 \leq i \leq M$ , with probability  $1 - \epsilon$ . Let  $\hat{r}_t := \sum_{j=1}^M q_{t,j} r_{t,j}$ . We may re-write the second term of the regret as follows,

$$\sum_{t=1}^{X^n} (1 - \tau_t) (r_{t,i} - r_t) = \sum_{t=1}^{X^n} (1 - \tau_t) (r_{t,i} - \hat{r}_t) + \sum_{t=1}^{X^n} (1 - \tau_t) (\hat{r}_t - r_t) + \sum_{t=1}^{X^n} (1 - \tau_t) (r_t - \hat{r}_t)$$

We bound the second term on the right hand side, using Lemma 17

$$\sum_{t=1}^{X^n} (1 - \tau_t) (\hat{r}_t - r_t) \leq \sum_{t=1}^{X^n} (\hat{r}_t - r_t) \frac{\log M}{n} + \sum_{t=1}^{X^n} \sum_{j=1}^{X^M} q_{t,j} r_{t,j}^2$$

As for the third term,

$$\sum_{t=1}^{X^n} (1 - \tau_t) \sum_{j=1}^{X^M} q_{t,j} r_{t,j} (r_t - r_{t,j}) = \sum_{t=1}^{X^n} (1 - \tau_t) \sum_{j=1}^{X^M} q_{t,j} (\hat{r}_t - r_{t,j} + r_{t,j} - r_t) = \sum_{t=1}^{X^n} (1 - \tau_t) \sum_{j=1}^{X^M} q_{t,j} \hat{r}_t - \sum_{t=1}^{X^n} (1 - \tau_t) \sum_{j=1}^{X^M} q_{t,j} r_{t,j} + \sum_{t=1}^{X^n} (1 - \tau_t) \sum_{j=1}^{X^M} q_{t,j} r_t$$

It remains to bound the deviation term. For  $t$  that satisfy  $\tau_t = 0$ , the action/model is selected according to  $\hat{r}_t$ , therefore the conditional expectation can be written as

$$E_{\tau_t=0} r_t = \sum_{j=1}^{X^M} q_{t,j} r_{t,j}$$

The sequence  $X_t := r_t - E_{t-1} r_t$  is a martingale difference sequence adapted to the history since for every  $t \geq 1$ ,

$$E_{t-1} X_t = E[r_t - E_{t-1} r_t | H_{t-1}] = 0:$$

Since  $r_t \in [B, 2B]$ , then  $X_t \in [-B, B]$  almost surely, which allows for an application of anytime Azuma-Hoeffding (Lemma 26):

$$P\left\{ \sum_{t=1}^n X_t \geq \frac{5B\sqrt{p}}{2} \sqrt{n((\log \log n B^2)_+ + \log(2B))} \right\} \leq \frac{1}{2}$$

which, in turn, leads us to

$$\sum_{t=1}^n (1 - \alpha_t) \sum_{j=1}^M \alpha_{t,j} r_{t,j} \leq \frac{5B\sqrt{p}}{2} \sqrt{n((\log \log n B^2)_+ + \log(2B))} \text{ w.h.p.}$$

simultaneously for all  $t \geq 1$ . We set  $\alpha_1 = \alpha_2 = \dots = \alpha_3$ , take a union bound and put the terms together obtaining,

$$R(n; i) \leq \sum_{t=1}^n \left( \frac{\log M}{n} + \sum_{j=1}^M \alpha_{t,j}^2 \right) + \frac{5B\sqrt{p}}{2} \sqrt{n((\log \log n B^2)_+ + \log(6B))} + 5B \sqrt{n((\log \log n)_+ + \log(12B))}$$

We upper bound the sum of last two terms to conclude the proof.  $\square$

The next two lemmas bound the bias and variance terms which appear in Lemma 18.

Lemma 19 (Anytime Bound on the Bias Term) If the regularization parameter of Lasso is chosen at every step as

$$\lambda_t = \frac{2}{\sqrt{t}} \left( 1 + \frac{5\sqrt{p}}{2} \sqrt{d(\log(2M) + (\log \log d)_+)} + \frac{12}{\sqrt{p}} (\log(2M) + (\log \log d)_+) \right)$$

and  $\lambda_t = O(t^{-1/4})$ , then with probability greater than  $1 - \frac{1}{n}$ , simultaneously for all  $t \geq 1$ ,

$$\sum_{t=1}^n \sum_{j=1}^M \alpha_{t,j}^2 r_{t,j} \leq n^{3/4} C_{\min}^{-1} C(M; d) \left( 1 + n^{-3/8} C_{\min}^{-1} \sqrt{p(\log(Md) + (\log \log n)_+)} \right)$$

where

$$C(M; d) := C \sqrt{1 + \frac{p}{d(\log(M) + (\log \log d)_+) + (\log(M) + (\log \log d)_+)}}$$

and  $C$  is an absolute constant.

Proof of Lemma 19 By the definition of the expected reward and its estimate,

$$\sum_{t=1}^n \sum_{j=1}^M \alpha_{t,j}^2 r_{t,j} = \sum_{t=1}^n \sum_{j=1}^M \alpha_{t,j}^2 \int r(x) \hat{\pi}_t(x) d p_{t+1,j}(x)$$

$$\stackrel{\text{c.s.}}{\leq} \sum_{t=1}^n \sum_{j=1}^M \alpha_{t,j}^2 \int r(x) \hat{\pi}_t(x) d p_{t+1,j}(x)$$

$$\stackrel{\text{bdd.}}{\leq} \sum_{t=1}^n \sum_{j=1}^M \alpha_{t,j}^2 \int \sum_{k=1}^K \hat{\pi}_t(x) k_2 d p_{t+1,j}(x) = \sum_{t=1}^n \sum_{j=1}^M \alpha_{t,j}^2 \int \sum_{k=1}^K d p_{t+1,j}(x)$$



Which implies for all  $t \geq 1$ ,

$$\sum_{j=1}^M \alpha_{t,j} r_{t,j}^2 \leq \frac{4^p \overline{10}^p}{c_{t,1}^2} + B \sum_{j=1}^M \alpha_{t,j} = \frac{160^p}{c_{t,1}^4} + B^2 + \frac{8B^p \overline{10}^p}{c_{t,1}^2}.$$

For the last term, similar to the proof of Lemma 19 we have,

$$\sum_{t=1}^n \frac{8B^p \overline{10}^p}{c_{t,1}^2} \leq C_1 B n^{1-4} + 1 + \frac{n^{3-8p}}{C_{\min}} \frac{1}{\log(Md) + (\log \log n)_+}$$

for some absolute constant  $C_1$ . We treat the squared term similarly,

$$\begin{aligned} \sum_{t=1}^n \frac{160^p}{c_{t,1}^4} &\leq C(M; ; d) \sum_{t=1}^n \frac{C_3}{t C_{\min}} + C_4 \frac{n^{3-8p}}{C_{\min}} \frac{1}{\log(Md) + (\log \log t)_+} \\ &\leq C(M; ; d) \frac{C_3 \log n}{C_{\min}} + 1 + C_4 \frac{n^{3-8p}}{C_{\min}} \frac{1}{\log(Md) + (\log \log n)_+} \end{aligned}$$

Note that the last inequality is not tight. This term will not be fastest growing term in the regret, so we have little motivation to bound it tightly. Therefore,

$$\begin{aligned} \sum_{t=1}^n \sum_{j=1}^M \alpha_{t,j} r_{t,j}^2 &\leq C_1 B^2 C_{\min}^p \overline{n} + C_2 B n^{1-4} + 1 + \frac{n^{3-8p}}{C_{\min}} \frac{1}{\log(Md) + (\log \log t)_+} \\ &\quad + C(M; ; d) \frac{\log n}{C_{\min}} + 1 + C_4 \frac{n^{3-8p}}{C_{\min}} \frac{1}{\log(Md) + (\log \log n)_+} \end{aligned}$$

where  $C_i$  are absolute constants. □

**Proof of Lemma 16** We start by conditioning on the event that  $t$  is picked such that  $\alpha_{t,j} \geq 1$  for all  $t \geq 1$  and  $j = 1, \dots, M$ . Then by application of Lemma 18 we get with probability greater than  $1 - 2^{-3}$ ,

$$\begin{aligned} R(n; i) &\leq 2B \sum_{t=1}^n \frac{\log M}{n} + \sum_{t=1}^n \sum_{j=1}^M \alpha_{t,j} r_{t,j}^2 + \sum_{t=1}^n \left( \sum_{j=1}^M \alpha_{t,j} \right) \\ &\quad + 10B^p \frac{1}{n} \left( (\log \log n B^2)_+ + \log(12) \right) \end{aligned}$$

We invoke Lemma 19 and Lemma 20 with  $\beta = 3$  take a union bound, to bound the variance and  $\sum_{t,j} \alpha_{t,j} r_{t,j}^2$  terms as well. These lemmas require one application of Theorem 3 to hold simultaneously and no additional union bound is required between them, since the randomness comes only from the confidence interval over  $\hat{\mu}_t$ .

$$\begin{aligned} R(n; i) &\leq C_1 B n^{3-4} + \frac{\log M}{n} \\ &\quad + C_2 B^2 C_{\min}^p \overline{n} + C_3 B n^{1-4} + 1 + \frac{n^{3-8p}}{C_{\min}} \frac{1}{\log(Md) + (\log \log n)_+} \\ &\quad + C(M; ; d) \frac{\log n}{C_{\min}} + 1 + \frac{n^{3-8p}}{C_{\min}} \frac{1}{\log(Md) + (\log \log n)_+} \\ &\quad + n^{3-4} C(M; ; d) + 1 + \frac{n^{3-8p}}{C_{\min}} \frac{1}{\log(Md) + (\log \log n)_+} \\ &\quad + 10B^p \frac{1}{n} \left( (\log \log n B^2)_+ + \log(12) \right) \end{aligned}$$

with probability greater than  $1 - 2^{-3}$ , conditioned on event  $\mathcal{E}$ . Assuming that event  $\mathcal{E}$  happens with probability  $1 - 2^{-3}$ , let  $B = B(\gamma; M; d; n; B; ; C_{\min})$  denote the right-hand-side of the regret inequality above.

By the chain rule we may write,

$$P \text{Reg}(n; i) \leq B$$

$$P \left( \sum_{t=1}^n R(n; i) \leq B \right) \geq 1 - P \left( \sum_{t=1}^n R(n; i) > B \right)$$

$$P \left( \sum_{t=1}^n R(n; i) \leq B \right) \geq 1 - P \left( \sum_{t=1}^n R(n; i) > B \right)$$

$$(1 - \epsilon)(1 - \epsilon) \geq 1 - 3\epsilon$$

It remains to verify that even  $E$  is met with probability  $1 - 2\epsilon$ . Recall that  $t = O(C_{\min} \sqrt{t} C(M; ; d))$ , and that from Lemma 8 with probability  $1 - \epsilon$ ,

$$\frac{C_{\min}}{4 \sqrt{t} C(M; ; d)} \leq \frac{C_1 C_{\min} t^{1/4} + C_2 t^{1/8} \sqrt{\log(Md) + (\log \log t)_+}}{4 \sqrt{10} C(M; ; d)} \leq \frac{2}{4 \sqrt{10} t}$$

Therefore, from Lemma 20, there exists  $t$  such that  $t = C C_{\min} \sqrt{t} C(M; ; d)$  satisfying,

$$P \left( \sum_{t=1}^n R(n; i) \leq B \right) \geq 1 - 2\epsilon$$

The proof is then finished by setting  $\epsilon = 3$  (and updating the absolute constants).  $\square$

## D.2 Proof of Virtual Regret

Proposition 21. For any fixed  $\epsilon > 0$ , there exists an absolute constant  $C_1$  such that

$$P \left( \sum_{t=1}^n R(n; i) \leq C_1 \sqrt{t} \right) \geq 1 - \epsilon$$

where

$$C_1 \sqrt{t} := C_1 \frac{r \sqrt{2d \log \frac{t}{d} + 2 \log(1 + B^2)}}{\sqrt{C_{\min} t^{3/4}}} + C_{\min} t^{3/8} \sqrt{\log(Md) + (\log \log t)_+}$$

Moreover, for  $u_{t,j}(\cdot) := \sum_{t=1}^n R(n; i) + \sum_{t=1}^n R(n; i)$ ,

$$P \left( \sum_{t=1}^n R(n; i) \leq C_1 \sqrt{t} \right) \geq 1 - \epsilon$$

Proof of Proposition 21. Define for convenience  $V_t = \sum_{t=1}^n R(n; i) + \sum_{t=1}^n R(n; i)$ . We first observe that

$$\sum_{t=1}^n R(n; i) = V_t^{-1} \left( \sum_{t=1}^n R(n; i) \right) \geq y_t$$

We can apply results from Abbasi-Yadkori et al. [2011] to get an anytime-valid confidence set. Their Theorem 2 asserts that with probability  $1 - \epsilon$ , for all  $t \geq 1$  we have

$$\sum_{t=1}^n R(n; i) \leq \frac{2}{V_t} \sqrt{t}$$

where

$$t = 2 \log \frac{\det(V_t)^{1/2}}{\det(I)^{1/2}} + B^2$$

Clearly,  $V_t \geq \min(V_t) I$ , and therefore with high probability,

$$\sum_{t=1}^n R(n; i) \leq \frac{2}{\min(V_t)} \sqrt{t}$$

uniformly over time. Our assumption is that  $\|j(x)\| \leq 1$ , and hence, denoting by the eigenvalues of  $V_t$ , the geometric-arithmetic mean inequality yields

$$\det(V_t) \geq \prod_{i=1}^d \lambda_i \geq \frac{1}{d} \text{trace}(V_t)^d$$

<sup>2</sup>Their theorem statement is slightly different, but they prove the stronger version we state below.

Given that

$$\text{trace}(V_t) = \sum_{i=1}^d \sum_{s=1}^t (j^s(x))_i^2 + \tau \quad t + \tau$$

we can conclude that

$$t \geq 2 \log \frac{(t+d-\tau)^{d-2}}{\tau^{d-2}} + B^2 = d \log \frac{t}{\tau} + 1 + 2 \log(1 + \frac{\tau}{t}) + B^2$$

We note that

$$\min(V_t) = \min(\sum_{t_j} t_j) + \tau$$

Then due to Lemma 9, there exist absolute constants  $C_1$  and  $C_2$  such that for all  $t \geq 1$ ,

$$\min(V_t) \geq \tau + C_1 C_{\min} t^{3-4} - C_2 t^{3-8} \frac{P}{\log(Md) + (\log \log t)}$$

therefore, there exist  $C_3$  and  $C_4$  such that

$$\begin{aligned} P \frac{1}{\min(V_t)} &\leq \frac{1}{\tau + C_1 C_{\min} t^{3-4} - C_2 t^{3-8} \frac{P}{\log(Md) + (\log \log t)}} \\ &\leq \frac{C_3}{\tau + C_1 C_{\min} t^{3-4}} \left( 1 + \frac{t^{3-8} \frac{P}{\log(Md) + (\log \log t)}}{\tau + C_1 C_{\min} t^{3-4}} \right) \\ &\leq \frac{C_4}{\tau + C_1 C_{\min} t^{3-4}} \left( 1 + C_{\min}^{-1} t^{3-8} \frac{P}{\log(Md) + (\log \log t)} \right) \end{aligned}$$

with high probability for all  $t \geq 1$ . Setting

$$! (t; ; d) = C_5 \frac{r}{\tau + C_{\min} t^{3-4}} \frac{2d \log(\frac{t}{\tau} + 1) + 2 \log(1 + \frac{\tau}{t}) + B^2}{\tau + C_{\min} t^{3-4}} \left( 1 + C_{\min}^{-1} t^{3-8} \frac{P}{\log(Md) + (\log \log t)} \right)$$

where  $C_5$  is an absolute constant concludes the parametric confidence bound. The upper confidence bound then simply follows: for any  $x \in \mathcal{X}$

$$r(x) \leq \sum_{t_j} \hat{r}_{t_j}(x) = \sum_{t_j} \hat{r}_{t_j}(x) \leq \sum_{t_j} \hat{r}_{t_j}(x) + \sum_{t_j} k_{j^s(x)} k_2 \leq ! (t; ; d)$$

where the last inequality holds with high probability simultaneously for all  $x$ .  $\square$

Proof of Lemma 15 Using Proposition 21 and the Cauchy-Schwarz inequality we obtain,

$$\begin{aligned} R_{j^s}(n) &= \sum_{t=1}^n r(x^t) - r(x_{t_j}) \\ &= \sum_{t=1}^n r(x^t) - \sum_{t=1}^n \hat{r}_t(x^t) + \sum_{t=1}^n \hat{r}_t(x^t) - \sum_{t=1}^n \hat{r}_t(x_{t_j}) + \sum_{t=1}^n \hat{r}_t(x_{t_j}) - r(x_{t_j}) \\ &\leq \sum_{t=1}^n \sum_{j^s} \hat{r}_{t_j} \left( k_{j^s(x^t)} k_2 + k_{j^s(x_{t_j})} k_2 \right) + \sum_{t=1}^n \hat{r}_t(x^t) - \sum_{t=1}^n \hat{r}_t(x_{t_j}) \\ &\leq ! (t; ; d) \sum_{j^s} \left( k_{j^s(x^t)} k_2 + k_{j^s(x_{t_j})} k_2 \right) + \sum_{t=1}^n \hat{r}_t(x^t) - \sum_{t=1}^n \hat{r}_t(x_{t_j}) \end{aligned}$$

with probability  $1 - \delta$ . If the agent selects actions greedily, then  $r(x^t) = \hat{r}_t(x_{t_j})$ , and

$$R_{j^s}(n) \leq \sum_{t=1}^n ! (t; ; d) \sum_{j^s} \left( k_{j^s(x^t)} k_2 + k_{j^s(x_{t_j})} k_2 \right) \leq 2 ! (t; ; d)$$

since the feature map is normalized to satisfy  $\sum_{j^s} k_{j^s} \leq 1$ . If the agent selects actions optimistically according to the upper confidence bound of Proposition 21, then

$$\hat{r}_t(x_{t_j}) + ! (t; ; d) \sum_{j^s} k_{j^s(x_{t_j})} k_2 \leq \hat{r}_t(x^t) + ! (t; ; d) \sum_{j^s} k_{j^s(x^t)} k_2$$

which implies

$$\mathbb{P}_t(x^?) \leq \mathbb{P}_t(x_{t;j}) \leq (t; d)k_j(x_{t;j}) \leq (t; d)k_j(x^?)k$$

and therefore,

$$\mathbb{R}_j^?(n) = \sum_{t=1}^n (t; d)k_j(x_{t;j})k_2 \leq \sum_{t=1}^n (t; d):$$

Then due to Proposition 21, with probability greater than  $1 - \epsilon$ , simultaneously for all  $j = 1, \dots, d$ ,

$$\sum_{t=1}^n (t; d) \leq C_1 \frac{n^{2d \log \frac{t}{d} + 1 + 2 \log(1 + \frac{1}{B^2})}}{C_{\min} t^{3=4}} + t^{3=8} \frac{q}{\log(\frac{Md}{t}) + (\log \log t)_+}$$

$$\leq C_1 n^{5=8} \frac{n^{2d \log \frac{n}{d} + 1 + 2 \log(1 + \frac{1}{B^2})}}{C_{\min}} + C_{\min}^{-1} n^{\frac{3}{8}q} \frac{q}{\log(\frac{Md}{n}) + (\log \log n)_+}$$

concluding the proof.  $\square$

## E Time-Uniform Concentration Inequalities

We will make use of the elegant concentration results in Howard et al. [2021], which analyzes the boundary of sub-Gamma processes.

**Definition 22 (Sub-Gamma process)** Let  $(S_t)_{t=0}^1$  and  $(V_t)_{t=0}^1$  be real-valued processes adapted to  $(\mathcal{F}_t)_{t=0}^1$  with  $S_0 = V_0 = 0$  and  $V_t$  non-negative. We say that  $(S_t)$  is sub-Gamma if for  $\epsilon \in [0, 1-c)$ , there exists a supermartingale  $(M_t)_{t=0}^1$  w.r.t.  $\mathcal{F}_t$ , such that  $\mathbb{E} M_0 = 1$  and for all  $t \in [0, 1]$ :

$$\exp(S_t - \frac{\epsilon}{2(1-c)} V_t) \leq M_t(\epsilon) \quad \text{a.s.}$$

The following is a special case of Theorem 1 in Howard et al. [2021]. We have simplified it by making a few straightforward choices for the parameters used originally by Howard et al. [2021], which will yield an easier-to-use bound in our scenario.

**Proposition 23 (Curved Boundary of Sub-Gamma Processes)** Let  $(S_t)_{t=0}^1$  be sub-Gamma with scale parameter  $c$  and variance process  $(V_t)_{t=0}^1$ . Define the boundary

$$B(\epsilon, v) := \frac{5}{2} \max\{v; 1\} g(\log \log v)_+ + \log \frac{2}{\epsilon} + 3c (\log \log v)_+ + \log \frac{2}{\epsilon};$$

for  $v > 0$ , where  $(x)_+ = \max(0; x)$ . Then,

$$\mathbb{P}(\exists t : S_t \geq B(\epsilon, V_t)) \leq \epsilon.$$

**Proof of Proposition 23.** Suppose  $\zeta(s)$  denotes the Riemann zeta function. Theorem 1 in Howard et al. [2021] states that  $(S_t)_{t=0}^1$  is a sub-Gamma process with variance process  $(V_t)_{t=0}^1$  then the boundary

$$S(v^0) = k_1 v^0 s \log \log(v^0) + \log \frac{(s)}{\log^s} + ck_2 s \log \log(v^0) + \log \frac{(s)}{\log^s} :$$

satisfies,

$$\mathbb{P}(\exists t : S_t \geq S(\max(V_t; 1))) \leq \epsilon,$$

where

$$k_1 := \frac{1=4}{p} + \frac{1=4}{2} \quad \text{and} \quad k_2 := (p^{-1} + 1) = 2$$

and; 1. Choosing  $s = 2$  and  $\epsilon = e$ , we obtain  $(2) = 2=6 = 2$ . Furthermore, we have  $k_1 \geq \frac{3}{2}$  and  $k_2 \geq \frac{3}{2}$ . Then if  $v^0 \geq 1$  (which we will enforce by the construction  $v^0 = \max(1; v)$ ), we compute

$$s \log \log(v^0) + \log \frac{(s)}{\log^s} \leq 2(\log \log v^0)_+ + \log \frac{2}{\epsilon} :$$



Therefore, we can upper bound (using our bounds  $k_1, k_2$ )

$$S(v^0) \leq \frac{5}{2} v^0 (\log \log v^0)_+ + \log \frac{2}{\epsilon} + 3c (\log \log v^0)_+ + \log \frac{2}{\epsilon} :$$

Now, since the boundary is given by  $S_t = (\max(v; 1))$  and  $v^0 = \max(v; 1) - 1$  we deduce that

$$B(v) := \frac{5}{2} \max(v; 1) (\log \log v)_+ + \log \frac{2}{\epsilon} + 3c (\log \log v)_+ + \log \frac{2}{\epsilon} :$$

is an any-time valid boundary.  $\square$

**Lemma 24 (Time-Uniform Two-sided Bernoulli)** Let  $X_1, \dots, X_s, \dots, X_t$  be a martingale sequence of Bernoulli random variables with conditional mean  $\mu_s$ . Then for all  $\epsilon > 0$ ,

$$P(9t : S_t \leq B(V_t)) \leq \frac{5\epsilon}{2} \frac{t}{t + (\log \log t)_+ + \log(4/\epsilon)} ;$$

**Proof of Lemma 24** By Proposition 23, we know that  $S_t$  is sub-Gamma with variance process  $V_t$  and scale parameter  $\epsilon$ , then

$$P(9t : S_t \leq B(V_t)) ;$$

where

$$B(v) := \frac{5\epsilon}{2} \max(1; v) (\log \log v)_+ + \log(2/\epsilon) + 3c (\log \log v)_+ + \log(2/\epsilon) :$$

By Howard et al. [2020], we know that  $(X_s)_{s=1}^t$  is a Bernoulli sequence, then  $S_t = \sum_{s=1}^t (X_s - \mu_s)$  is sub-Gamma with variance process  $V_t = t$  and scale parameter  $\epsilon = 0$  (hence, sub-Gaussian). This implies,

$$P(9t : S_t \leq B(V_t)) \leq \frac{5\epsilon}{2} \frac{t}{t + (\log \log t)_+ + \log(2/\epsilon)} ;$$

The above arguments also holds for the sequence  $X_s$ . Then taking a union bound and adjusting  $\epsilon = 2\epsilon$  concludes the proof.  $\square$

**Lemma 25 (Time-Uniform Bernstein)** Let  $(\epsilon_i)_{i=1}^t$  be a sequence of conditionally standard sub-gaussian variables, where each  $\epsilon_i$  is  $F_{i-1} = (\epsilon_1, \dots, \epsilon_{i-1})$  measurable. Then, for  $\epsilon_i \in \mathbb{R}$  and  $v_i \in (0, 1]$

$$P(9t : \sum_{i=1}^t (\epsilon_i^2 - v_i) \leq \frac{5}{2} \max_{i=1, \dots, t} (4kv_t k_2^2 + (\log \log t)_+ + \log(2/\epsilon)) \max_{i=1, \dots, t} v_i) \leq \epsilon$$

where,  $v_t = (v_1, \dots, v_t) \in \mathbb{R}^t$  and  $! (v) := \log \log(4\epsilon v^2)_+ + \log(2/\epsilon)$ .

**Proof of Lemma 25** From Lemma 4,  $S_t = \sum_{i=1}^t (\epsilon_i^2 - v_i)$  is sub-Gamma with variance process  $V_t = 4 \sum_{i=1}^t v_i^2$  and  $c = 4 \max_{i=1, \dots, t} v_i$ . By Proposition 23, we know that  $S_t$  is sub-Gamma with variance process  $V_t$  and scale parameter  $\epsilon$ , then

$$P(9t : S_t \leq B(V_t)) ;$$

where

$$B(v) := \frac{5\epsilon}{2} \max(1; v) (\log \log v)_+ + \log(2/\epsilon) + 3c (\log \log v)_+ + \log(2/\epsilon) :$$

$\square$

**Lemma 26 (Time-Uniform Azuma-Hoeffding)** Let  $X_1, \dots, X_n$  be a martingale difference sequence such that  $|X_t| \leq B$  for all  $t > 1$  almost surely. Then for all  $\epsilon > 0$ ,

$$P(9t : \sum_{s=1}^t X_s \leq \frac{5B}{2} \frac{t}{t + (\log \log t)_+ + \log(2/\epsilon)}) \leq \epsilon ;$$

*Proof of Lemma 26.* By Proposition 23, we know that if  $S_t$  is sub-Gamma with variance process  $V_t$  and scale parameter  $c$ , then

$$\mathbb{P}(S_t \geq B_\delta(V_t)) \leq \delta,$$

where

$$B_\delta(v) := \frac{5}{2} \sqrt{\max\{1, v\} (\log \log ev)_+ + \log(2/\delta)} + 3c (\log \log ev)_+ + \log(2/\delta).$$

By [Howard et al., 2020], we know that if  $(X_t)_{t=1}^T$  is  $B$ -bounded martingale difference sequence, then  $S_t = \sum_{s=1}^t X_s$  is sub-Gamma with variance process  $V_t = tB^2$  and scale parameter  $c = 0$ . This implies,

$$\mathbb{P}(S_t \geq \frac{5B}{2} \sqrt{t((\log \log etB^2)_+ + \log(2/\delta))}) \leq \delta,$$

concluding the proof.  $\square$

## F Experiment Details

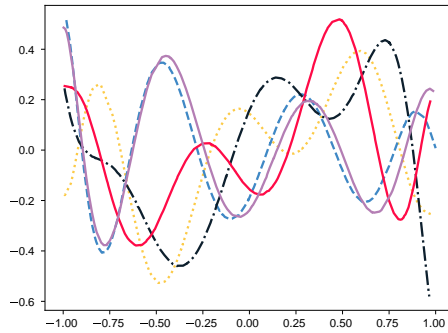


Figure 5: Examples of possible reward functions  $r(\cdot)$  in our experiments.

### F.1 Hyper-Parameter Tuning Results

We implement 6 algorithms in our experiments, ETC [Algorithm 4, Hao et al., 2020], ETS (Algorithm 5), CORRAL [Algorithm 6, Agarwal et al., 2017], ALEXP (Algorithm 1), and Lastly UCB (Algorithm 3) with the oracle feature map  $j_*$  (Oracle), and UCB with the concatenated feature map (Naive). The Python code is available on [github.com/lasgroup/ALEXP](https://github.com/lasgroup/ALEXP). When algorithms require exploration, e.g., in the case of ETC or ALEXP, we simply set  $\pi = \text{Unif}(X)$ . Figure 7 shows the results of our hyperparameter tuning experiment. To ensure that the curves are valid, we run each

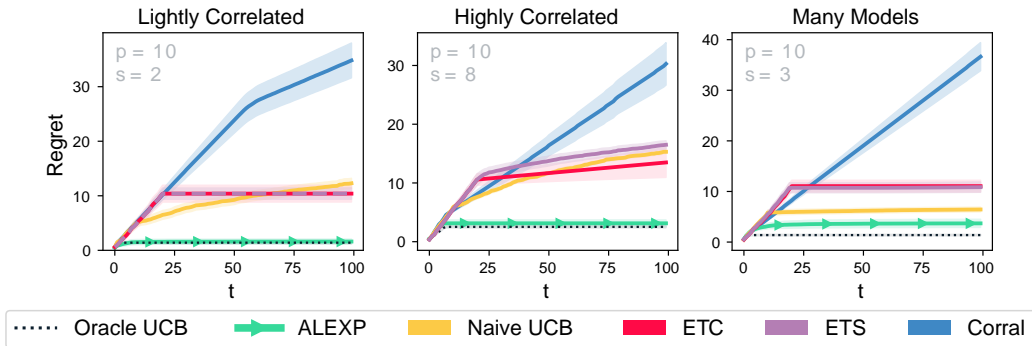


Figure 6: Bench-marking ALEXP and other baselines. Complete version of Fig. 1 and Fig. 2.

---

**Algorithm 2** GetPosterior

---

Inputs:  $H_t, \tilde{\lambda}$   
 Let  $K_t = [ \sum_{i,j=1}^p \langle \mathbf{x}_i, \mathbf{x}_j \rangle ]_{i,j=1}^p$ , and  $V_t = (K_t + \tilde{\lambda}^2 I)$ , and  $\mathbf{k}(\cdot) = [ \sum_{i=1}^p \langle \mathbf{x}_i, \cdot \rangle ]_{i=1}^p$

Calculate  $\mu_t(\cdot) = \frac{\mathbf{k}^T(\cdot) V_t^{-1} \mathbf{y}_t}{\mathbf{1}^T V_t^{-1} \mathbf{1}}$

Calculate  $\sigma_t(\cdot) = \sqrt{\mathbf{1}^T V_t^{-1} \mathbf{k}(\cdot) \mathbf{k}^T(\cdot) V_t^{-1} \mathbf{k}(\cdot)}$

Return:  $\mu_t, \sigma_t$

---



---

**Algorithm 3** UCB

---

Inputs:  $\tilde{\lambda}, \beta_t$ ,  
**for**  $t = 1, \dots, n$  **do**  
 Choose  $\mathbf{x}_t = \arg \max_{\mathbf{x}} u_{t-1}(\mathbf{x}) = \mu_{t-1}(\mathbf{x}) + \beta_t \sigma_{t-1}(\mathbf{x})$ . ▷ Choose actions optimistically  
 Observe  $y_t = r(\mathbf{x}_t) + \varepsilon_t$ . ▷ Receive reward  
 $H_t = H_{t-1} + [ \sum_{i=1}^p \langle \mathbf{x}_i, y_t \rangle ]_{i,j=1}^p$  ▷ Append history  
 Update  $\mu_t, \sigma_t = \text{GetPosterior}(H_t, \tilde{\lambda})$   
**end for**

---

configuration for 20 different random seeds, i.e. on different random environments. The shaded areas in Figure 7 show the standard error.

**UCB.** For all the experiments, we set the exploration coefficient of UCB to  $\beta_t = 2^3$  and choose the regression regularizer from  $\tilde{\lambda} \in \{0.01, 0.1, 0.5\}g$ . We use PYTORCH [Paszke et al., 2017] for updating the upper confidence bounds, which requires more regularization for longer feature maps (e.g. when  $s = 8, p = 2$ ), to be computationally stable.

**Lasso.** Every time we need to solve Eq. (1), we set  $\lambda_t$  according to the *rate* suggested by Theorem 3. To find a suitable constant scaling coefficient, we perform a hyper-parameter tuning experiment sampling 20 values in  $[10^{-5}, 10^0]$ . We choose  $\lambda_0 = 0.009$ , and scale  $\lambda_t$  with it across all experiments.

**ALEXP.** We set the rates for  $\gamma_t$  and  $\eta_t$  as prescribed by Theorem 1. For the scaling constants, we perform a hyper-parameter tuning experiment log-uniformly sampling 20 different configurations from  $\gamma_0 \in [10^{-4}, 10^{-1}]$  and  $\eta_0 \in [10^0, 10^2]$ . For each problem instance (i.e. as  $s$  and  $p$  change) we repeat this process. However we observe that the optimal hyper-parameters work well across all problem instances.

**ETC/ETS.** For these algorithms, we separately tune  $n_0$  for each problem instance. We set  $\lambda_1 \propto \frac{1}{\log M/n_0}$  according to Theorem 4.2 of [Hao et al., 2020] and scale it with  $\lambda_0 = 0.009$ , as stated before. We uniformly sample 10 different values where  $n_0 \in [2, 80]$  since the horizon is  $n = 100$ . The optimal value often happens around  $n_0 = 20$ .

**CORRAL.** We set the rates of the parameters as  $\gamma = O(1/n)$  and  $\eta = O(\frac{1}{\sqrt{M/n}})$  according to Agarwal et al. [2017, Theorem 5,]. Then similar to ALEXP, we tune the scaling constants. The procedure for tuning the constants is identical to ALEXP, as in we use the same search interval, and try 10 different configurations for  $\gamma$  and  $\eta$ .

---

<sup>3</sup>To achieve the  $\frac{P}{dT \log T}$  regret, one has to set  $\lambda_t = O(\frac{P}{d \log T})$  as shown in Proposition 21.

---

**Algorithm 4** ETC [Hao et al., 2020]

---

Inputs:  $n_0, n, \lambda_1, \pi$   
Let  $H_0 = ;$   
**for**  $t = 1, \dots, n_0$  **do**  
    Draw  $\mathbf{x}_t \sim \pi$ . ▷ Explore.  
    Observe  $y_t = r(\mathbf{x}_t) + \varepsilon_t$ . ▷ Receive reward  
     $H_t = H_{t-1} \cup [f(\mathbf{x}_t, y_t)g]$  ▷ Append history  
**end for**  
 $\hat{L}_{n_0} = L(\cdot, H_{n_0}, \lambda_1)$  ▷ Perform Lasso once  
**for**  $t = n_0 + 1, \dots, n$  **do**  
    Choose  $\mathbf{x}_t = \arg \max_{\mathbf{x}} \hat{L}_{n_0}(\mathbf{x})$  ▷ Choose actions greedily  
**end for**

---

---

**Algorithm 5** ETS

---

Inputs:  $n_0, n, \lambda_1, \tilde{\lambda}, \beta_t, \pi$   
Let  $H_0 = ;$   
**for**  $t = 1, \dots, n_0$  **do**  
    Draw  $\mathbf{x}_t \sim \pi$ . ▷ Explore  
    Observe  $y_t = r(\mathbf{x}_t) + \varepsilon_t$ . ▷ Receive reward  
     $H_t = H_{t-1} \cup [f(\mathbf{x}_t, y_t)g]$  ▷ Append history  
**end for**  
 $\hat{L}_{n_0} = L(\cdot, H_{n_0}, \lambda_1)$  ▷ Perform Lasso once  
 $\hat{J} = \{j \mid \hat{L}_{n_0, j} \notin \mathbf{0}, j \in [M]g\}$  ▷ Get sparsity pattern  
 $\hat{J}(\cdot) = [\hat{L}_{n_0, j}(\cdot)]_{j \in \hat{J}}$  ▷ Model-select acc. to  $\hat{J}$   
**for**  $t = n_0 + 1, \dots, n$  **do**  
    Choose  $\mathbf{x}_t = \arg \max_{\mathbf{x}} u_{t-1}(\mathbf{x}) = \mu_{t-1}(\mathbf{x}) + \beta_t \sigma_{t-1}(\mathbf{x})$  ▷ Choose actions optimistically  
    Observe  $y_t = r(\mathbf{x}_t) + \varepsilon_t$   
     $H_t = H_{t-1} \cup [f(\mathbf{x}_t, y_t)g]$   
    Update  $\mu_t, \sigma_t = \text{GetPosterior}(H_t, \hat{J}, \tilde{\lambda})$   
**end for**

---

---

**Algorithm 6** CORRAL [Agarwal et al., 2017]

---

Inputs:  $n, \gamma, \eta$   
Initialize  $\beta = e^{1/\ln n}$ ,  $\eta_{1,j} = \eta$ ,  $\rho_{1,j} = 2M$  for all  $j = 1, \dots, M$   
Set  $\mathbf{q}_1 = \bar{\mathbf{q}}_1 = \frac{1}{M}$  and initialize base agents  $(p_{1,1}, \dots, p_{1,M})$ .  
**for**  $t = 1, \dots, n$  **do**  
    Choose  $j_t \sim \bar{\mathbf{q}}_t$ . ▷ Sample Agent  
    Draw  $\mathbf{x}_t \sim p_{t,j_t}$ . ▷ Play action according to agent  $j_t$   
    Observe  $y_t = r(\mathbf{x}_t) + \varepsilon_t$ .  
    Calculate IW estimates  $\hat{r}_{t,j} = \frac{y_t}{\bar{q}_{t,j}} \mathbb{1}_{f_j = j_t} g$  for all  $j = 1, \dots, M$ .  
    Send  $\hat{r}_{t,j} = \frac{y_t}{\bar{q}_{t,j}} \mathbb{1}_{f_j = j_t} g$  to agents and get updated policies  $p_{t+1,j}$ .  
     $\mathbf{q}_{t+1} = \text{LOG-BARRIER-OMD}(\mathbf{q}_t, \hat{r}_{t,j_t} \mathbf{e}_{j_t}, t)$  ▷ Update agent probabilities  
     $\bar{\mathbf{q}}_{t+1} = (1 - \gamma) \bar{\mathbf{q}}_{t+1} + \gamma \frac{1}{M}$  ▷ Mix with exploratory distribution  
    **for**  $j = 1, \dots, M$  **do** ▷ Update parameters  
        **if**  $\frac{1}{\bar{q}_{t+1,j}} > \rho_{t,j}$  **then**  $\rho_{t+1,j} = \frac{2}{\bar{q}_{t,j}}$ , and  $\eta_{t+1,j} = \beta \eta_{t,j}$   
        **else**  $\rho_{t+1,j} = \rho_{t,j}$  and  $\eta_{t+1,j} = \eta_{t,j}$   
        **end if**  
    **end for**  
**end for**

---

---

**Algorithm 7** LOG-BARRIER-OMD

---

Inputs:  $\mathbf{q}_t, \eta_t, t$   
Find  $\xi \in [\min_j \ell_{t,j}, \max_j \ell_{t,j}]$  such that  $\sum_{j=1}^M q_{t,j}^1 + \eta_{t,j}(\ell_{t,j} - \xi)^1 = 1$   
Return:  $\mathbf{q}_{t+1}$  where  $q_{t+1,j}^1 = q_{t,j}^1 + \eta_{t,j}(\ell_{t,j} - \xi)$  for all  $j \in [M]$

---

