



Zhou Lu<sup>1</sup>, Richard Zhang<sup>2</sup>, Xinyi Chen<sup>1,2</sup>, Fred Zhang<sup>3</sup>, David Woodruff<sup>2,4</sup>, Elad Hazan<sup>1,2</sup>  
 {<sup>1</sup> Princeton University, <sup>2</sup> Google, <sup>3</sup> UC Berkeley, <sup>4</sup> Carnegie Mellon University}

Introduction

## Background/Related Works

Motivation:

- Fast changing states pose a significant challenge to online optimization
- Want to perform rapid adaptation under **limited observation**
- The classic metric of regret incentivizes static behavior and is not correct in changing environments.
- Previous works proposed the notion of (strongly) adaptive regret, defined as the maximum regret over any continuous interval in time.

Can we efficiently learn the best learning rate/optimizer that adapts to changing environments (i.e. learn a schedule)?

$$\text{SA-regret}(\mathcal{A}, I) = \max_{s-j=I} \left[ \sum_{t=j}^s \ell_t^\top e_{x_t} - \min_i \sum_{t=j}^s \ell_t^\top e_i \right]$$

Table 1: Adaptive regret bounds and query efficiency in the adversarial multi-armed bandits setting.

Algorithm	Adaptive regret bound	Number of queries
FLH Hazan & Seshadhri (2009)	$\sqrt{nT}$	$O(\log T)$
SAOL Daniely et al. (2015)	$\sqrt{nI \log T}$	$O(\log T)$
EFLH Lu & Hazan (2023)	$I^{\frac{1}{2}+\epsilon} \cdot \sqrt{n \log T}$	$O\left(\frac{\log \log T}{\epsilon}\right)$
<b>This paper (Theorem 1)</b>	$\sqrt{nI \log n} \cdot \log^{1.5} T$	2

Method

## Main Algorithm (StABL)

**Algorithm 1** Strongly Adaptive Bandit Learner (StABL)

- 1: **Input:** general EXP3 algorithm  $\mathcal{A}$  and horizon  $T$ .
- 2: Construct geometric interval set  $S = \{[s2^k, (s+1)2^k - 1] \mid 0 \leq k \leq \log T, s \in \mathbb{N}^+\}$ .
- 3: Construct  $B = \log T$  independent instances of EXP3 algorithm  $\mathcal{A}_k$ , where  $\mathcal{A}_k$  optimizes each  $\{I \in S \mid 2^k = |I|\}$  one after another since they don't overlap.
- 4: Denote  $w_t(k)$  to be the weight assigned to  $\mathcal{A}_k$  at time  $t$  by the meta-algorithm.
- 5: Denote  $v(t, k) \in \mathbb{R}^n$  to be the distribution over arms by  $\mathcal{A}_k$  at time  $t$ , and  $v(t, k)_i$  to be the probability of sampling arm  $i$  by  $\mathcal{A}_k$  at time  $t$ .
- 6: Define  $\eta_k = \min\{1/2\sqrt{n}, 1/\sqrt{n|2^k|}\}$ , and initialize  $w_1(k) = \eta_k$  for all  $k \in [B]$ .
- 7: **for**  $\tau = 1, \dots, T$  **do**
- 8:   Let  $W_t = \sum_k w_t(k)$  and  $p(t) = \frac{1}{W_t}(\dots, w_t(k), \dots)$  be the distribution over the base learners.
- 9:   For all  $i \in [n]$ , let
 

$$P(t)_i = \frac{\max_k v(t, k)_i^2}{2 \sum_i \max_k v(t, k)_i^2} + \frac{\sum_k v(t, k)_i}{2B}$$
- 10:   // Observe that this defines a probability distribution over  $n$  arms.  
 Sample  $x_t \sim \sum_k p(t)_k v(t, k)$ , and in parallel sample  $x'_t \sim P(t)$ .  
 // Only the second sample  $x'_t$  will be used for weight updating.
- 11:   Play  $x_t$ , suffer loss  $\ell_t^\top e_{x_t}$  and observe loss  $\ell_t^\top e_{x'_t}$ . Compute loss estimator
 

$$\hat{\ell}_t = \mathbf{1}_{i=x'_t} \frac{1}{P(t)_{x'_t}} \ell_t^\top e_{x'_t}$$
- 12:   Update each EXP3 instance with the loss estimator  $\hat{\ell}_t$ , via Algorithm 2.
- 13:   Update the meta-algorithm's weights over base learners via the loss estimator  $\hat{\ell}_t$ . For each  $k$ , update  $w_{t+1}(k)$  as follows,
 

$$w_{t+1}(k) = \begin{cases} \eta_k & 2^k |t+1 \\ w_t(k) (1 + \eta_k \tilde{r}_t(k)) & \text{else} \end{cases}$$

where  $\tilde{r}_t(k) = \hat{\ell}_t^\top \sum_k p(t)_k v(t, k) - \hat{\ell}_t^\top v(t, k)$ .
- 14: **end for**

### Algorithm Overview

- EXP3-type algorithms for both the black-box base learners and the meta-algorithm
- Directly using EXP3 in the MAB setting will fail, because the weight distribution might become unbalanced
- Use addition observation to create unbiased estimators of the loss vector with controllable variance
- A naïve (but sub-optimal) choice is to sample uniformly over the arms with  $O(n^2)$  variance
- Use importance sampling to reduce variance to  $O(n)$

Results

## Two Queries Suffice

- In the bandit setting, there is a  $\Omega(\log T)$  lower bound for SA-regret!
  - Consider two arms: each interval with sublinear regret guarantee requires trying both arms.
- However, with two queries per round, we can achieve an optimal  $\sqrt{\log T}$  bandit algorithm.
- Crucially decouples the action and observation distribution, unlike EXP4 or previous algorithms
- For multi-arm bandit, we achieve optimal dependence on the number of arms.

**Theorem 1** (Adaptive regret minimization for multi-armed bandits). *For the multi-armed bandits problem with  $n$  arms and  $T$  rounds, Algorithm 1 achieves an expected adaptive regret bound of  $O(\sqrt{nI \log n} \log^{1.5} T)$ , using two queries per round.*

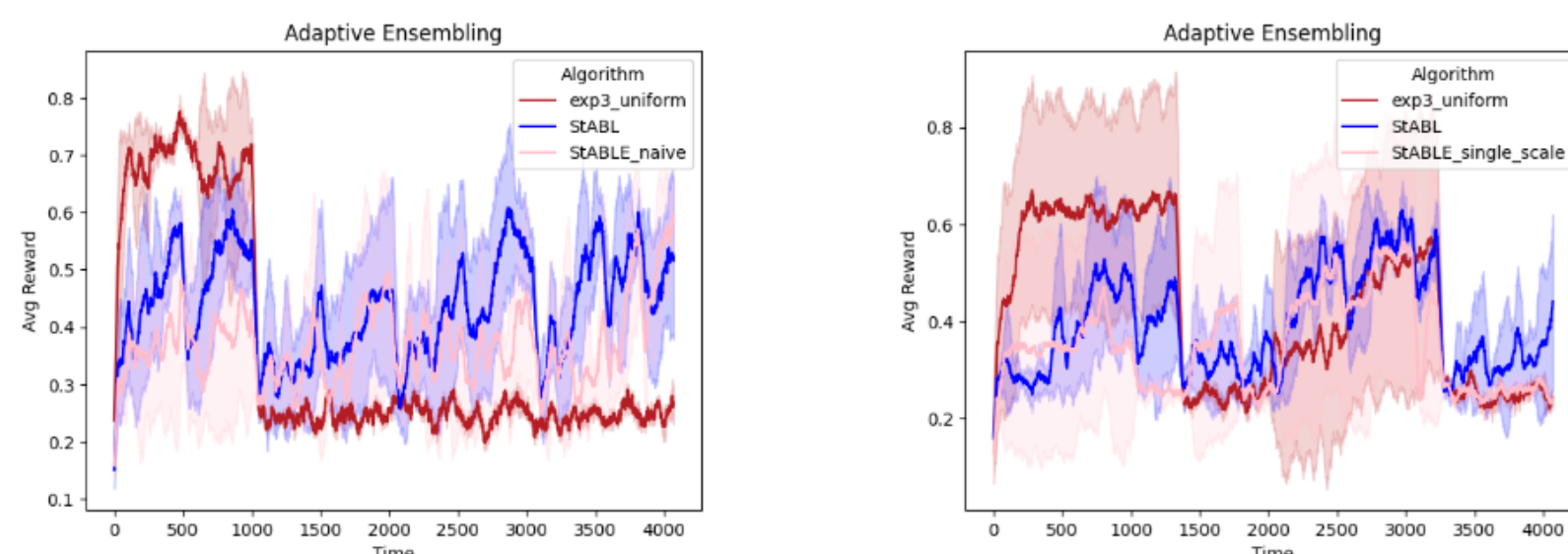


Figure 1: Comparison plots of the algorithm rewards in the learning with expert advice setting.

Conclusion

## Conclusion

We study adaptive regret under the limited observation model. Our result not only

- improves the state-of-the-art query efficiency of  $O(\log \log T)$  in Lu & Hazan (2023)
- matches the lower bound of the bandit setting Daniely et al. (2015)
- providing a sharp characterization of the query efficiency of adaptive regret.

This has multiple implications for:

- Learning rate adaptation
- Learning to optimize
- Meta-learning
- Adaptive gradient methods
- Data-driven algorithm design

## Next Steps

- **Explore vs Exploit:** UCB coefficient, random exploration/restarts, algorithm selection between explorative vs exploitative.
- **Multimetric Optimization:** Scalarizations can adaptively explore the Pareto frontier, especially when parts of the Pareto frontier are targeted.
- **Transfer Learning:** Balancing between the algorithms that 1) uses all the prior data and algorithms that 2) ignores all the prior data
- **Batch Setting:** Ensembling acquisitions and strategies in parallel in data-starved environments.
- **Early Stopping:** Early stopping is trading off between 1) never stopping for maximum exploration and 2) stopping early to save resources.