

# Non-stationary Bandit Convex Optimization: A Comprehensive Study

Xiaoqi Liu\*   Dorian Baudry\*   Julian Zimmert†   Patrick Rebeschini\*   Arya Akhavan\*

## Non-stationary Bandit Convex Optimization

- Adversary fixes convex loss functions  $f_1, f_2, \dots, f_T : \mathbb{R}^d \rightarrow [-1, 1]$
- For  $t \geq 1$ , learner
  - Selects action  $\mathbf{z}_t$  from continuous (convex & compact) arm set  $\Theta \subseteq \mathbb{R}^d$
  - Incurs loss  $f_t(\mathbf{z}_t)$ , observes *bandit feedback* with sub-Gaussian noise  $\xi_t$

$$y_t = f_t(\mathbf{z}_t) + \xi_t$$

- Regret* as benchmark with comparators  $\mathbf{u}_{1:T} = \{\mathbf{u}_1, \dots, \mathbf{u}_T\}$ :

$$R(T, \mathbf{u}_{1:T}) := \sum_{t=1}^T \mathbb{E} [f_t(\mathbf{z}_t) - f_t(\mathbf{u}_t)]$$

### Non-stationarity measures:

- Number of switches:  $S(\mathbf{u}_{1:T}) := 1 + \sum_{t=2}^T \mathbf{1}\{\mathbf{u}_t \neq \mathbf{u}_{t-1}\} \leq S$
- Path-length:  $P(\mathbf{u}_{1:T}) := \sum_{t=2}^T \|\mathbf{u}_t - \mathbf{u}_{t-1}\| \leq P$
- Total variation:  $\Delta(f_{1:T}) := \sum_{t=2}^T \max_{\mathbf{z} \in \Theta} |f_t(\mathbf{z}) - f_{t-1}(\mathbf{z})| \leq \Delta$

**Regret notions:** **switching**, **path-length**, and **dynamic regret**:

$$R^{\text{swi}}(T, S) := \max_{\mathbf{u}_{1:T}: S(\mathbf{u}_{1:T}) \leq S} R(T, \mathbf{u}_{1:T}), \quad R^{\text{path}}(T, P) := \max_{\mathbf{u}_{1:T}: P(\mathbf{u}_{1:T}) \leq P} R(T, \mathbf{u}_{1:T}),$$

$$R^{\text{dyn}}(T, \Delta) := \sup_{f_{1:T}: \Delta(f_{1:T}) \leq \Delta} \sum_{t=1}^T \mathbb{E} \left[ f_t(\mathbf{z}_t) - \min_{\mathbf{z} \in \Theta} f_t(\mathbf{z}) \right].$$

## Our goals

- Unified treatment** for non-stationary BCO (previous work: only  $R^{\text{dyn}}(T, \Delta)$  [1, 2] or only  $R^{\text{path}}(T, P)$  [3, 4])
- Design algorithms with optimal sublinear regret w.r.t.  $T$ ,  $S$ ,  $\Delta$  and  $P$

## Main results

Regret bounds for  $R^{\text{swi}}(T, S)$ ,  $R^{\text{dyn}}(T, \Delta)$  and  $R^{\text{path}}(T, P)$ , respectively, for algorithms *tuned with known*  $S$ ,  $\Delta$  and  $P$ :

Algo.:	TEWA-SE	cExO
General convex (GC)	$\sqrt{d}S^{\frac{1}{4}}T^{\frac{3}{4}}, d^{\frac{2}{5}}\Delta^{\frac{1}{5}}T^{\frac{4}{5}}, d^{\frac{2}{5}}P^{\frac{1}{5}}T^{\frac{4}{5}}$	
Strongly convex (SC)	$d\sqrt{ST}, d^{\frac{2}{3}}\Delta^{\frac{1}{3}}T^{\frac{2}{3}}, d^{\frac{2}{3}}P^{\frac{1}{3}}T^{\frac{2}{3}}$	$d^{\frac{5}{2}}\sqrt{ST}, d^{\frac{5}{3}}\Delta^{\frac{1}{3}}T^{\frac{2}{3}}, d^{\frac{5}{3}}\underline{P^{\frac{1}{3}}T^{\frac{2}{3}}}$
Comp. complexity	polynomial	exponential

- Straight underline: **minimax-optimal** rates.
- Wavy underline: result is either new to the literature (SC case) or **improves on the best-known**  $P^{\frac{1}{4}}T^{\frac{3}{4}}$  rate [3] (GC case).

## Algorithm 1: TEWA-SE

Tilted Exponentially Weighted Average with Sleeping Experts

**Input:** perturbation step-size  $h = \sqrt{d}B^{-\frac{1}{4}}$ , expert algorithm  $E(l, \eta)$  is *online gradient descent* over interval  $l$  with step-size parameterized by  $\eta$

```

1: for  $t = 1, 2, \dots, T$  do
2:   for Active expert  $E_i \equiv E_i(l_i, \eta_i) \in \{E_1, E_2, \dots, E_{n_t}\}$  do
3:     Receive action  $\mathbf{x}_{t,l_i}^{\eta_i}$  from expert  $E_i$ 
4:   end for
5:   Set meta-action using TEWA:
      
$$\mathbf{x}_t = \sum_{i=1}^{n_t} \frac{\eta_i e^{-L_{t-1,l_i}^{\eta_i}}}{\sum_{j=1}^{n_t} \eta_j e^{-L_{t-1,l_j}^{\eta_j}}} \mathbf{x}_{t,l_i}^{\eta_i} \quad \triangleright \text{aggregation}$$

6:   Sample  $\zeta_t$  uniformly from unit sphere  $\partial\mathbb{B}^d$ 
7:   Query point  $\mathbf{z}_t = \mathbf{x}_t + h\zeta_t$  to obtain  $y_t = f_t(\mathbf{z}_t) + \xi_t$ 
8:   Construct gradient estimate  $\mathbf{g}_t = (d/h)y_t\zeta_t$ 
9:   for  $i = 1, 2, \dots, n_t$  do
10:    Send meta-action  $\mathbf{x}_t$  and  $\mathbf{g}_t$  to  $E_i$ 
11:    Increment loss  $L_{t,l_i}^{\eta_i} = L_{t-1,l_i}^{\eta_i} + \ell_t^{\eta_i}(\mathbf{x}_{t,l_i}^{\eta_i}) \quad \triangleright \text{update experts}$ 
12:   end for
13: end for

```

**Our contribution: Construct SC surrogate losses with one-point gradient estimates:**

$\therefore$  Simple linear surrogate loss  $\ell_t(\mathbf{x}) = -\mathbf{g}_t^\top(\mathbf{x}_t - \mathbf{x}) \Rightarrow \sqrt{|I|}$  expert static regret  $\Rightarrow$  linear  $R^{\text{ada}}(B, T)$ .

$\therefore$  We instead use the SC surrogate loss

$$\ell_t^{\eta}(\mathbf{x}) = -\eta \mathbf{g}_t^\top(\mathbf{x}_t - \mathbf{x}) + \eta^2 G^2 \|\mathbf{x}_t - \mathbf{x}\|^2, \quad \forall \mathbf{x} \in \mathbb{R}^d$$

where w.h.p.  $\|\mathbf{g}_t\| \leq G \ \forall t \in [T] \Rightarrow \log |I|$  expert static regret  $\Rightarrow$  sublinear  $R^{\text{ada}}(B, T)$ .

Gradient estimate  $\mathbf{g}_t$  satisfies  $\mathbb{E}[\mathbf{g}_t | \mathbf{x}_t] = \nabla \hat{f}_t(\mathbf{x}_t)$  where  $\hat{f}_t(\mathbf{x}) = \mathbb{E}[f_t(\mathbf{x} + h\tilde{\zeta})]$  is smoothed loss ( $\tilde{\zeta}$  uniform on unit ball  $\mathbb{B}^d$ )  
Handle **bias-variance tradeoff** in regret upper bound by tuning  $h$

### Tools from prior work:

- Sleeping experts** on geometric intervals with geometric step-size
- TEWA aggregation to adapt to **unknown loss curvature** [5–7]

**Theorem:** For known  $B$ ,  $R^{\text{ada}}(B, T) \lesssim \begin{cases} \sqrt{d}B^{\frac{3}{4}} & (\text{GC}) \\ \frac{d}{\alpha}\sqrt{B} & (\text{SC}) \end{cases}$

## Minimax-optimal lower bounds

$$\overline{d\sqrt{ST} \quad d^{\frac{2}{3}}\Delta^{\frac{1}{3}}T^{\frac{2}{3}} \quad d^{\frac{4}{5}}P^{\frac{2}{5}}T^{\frac{3}{5}}}$$

- For  $d = 1$ , rates w.r.t.  $T, S, \Delta$  match those for multi-armed bandits
- Path-length bound improves on the only existing  $d\sqrt{PT}$  from [3]

## Algorithm 2: cExO clipped Exploration by Optimization

Vanilla ExO from [8] + **clipping**.

**Input:** a finite covering set  $\mathcal{C} \subset \Theta$  of  $\Theta$ , and  $\tilde{\Delta} = \Delta(\mathcal{C}) \cap [\gamma, 1]^{|\mathcal{C}|}$

```

1: for  $t = 1, \dots, T$  do
2:   Compute reference dist.:
      
$$\mathbf{q}_t = \Pi_{\tilde{\Delta}}(\tilde{\mathbf{q}}_t), \quad \tilde{\mathbf{q}}_t(\mathbf{x}) = \frac{e^{-\eta L_{t-1}(\mathbf{x})}}{\sum_{\mathbf{x}' \in \mathcal{C}} e^{-\eta L_{t-1}(\mathbf{x}')}} \quad \forall \mathbf{x} \in \mathcal{C}.$$

3:   Select sampling dist.  $\mathbf{p}_t \in \Delta(\mathcal{C})$  & loss estimator  $E_t$  by solving
      
$$\arg \min_{\substack{\mathbf{p} \in \Delta(\mathcal{C}), \\ E: \mathcal{C} \times [-1, 1] \rightarrow \mathbb{R}^{|\mathcal{C}|}}} \Lambda(\mathbf{q}_t, \mathbf{p}, E) \quad \triangleright \text{ExO step}$$

4:   Sample  $\mathbf{z}_t \sim \mathbf{p}_t$ , observe  $f_t(\mathbf{z}_t)$ 
5:   Set  $\ell_t = E_t(\mathbf{z}_t, f_t(\mathbf{z}_t))$ ,  $L_t(\mathbf{x}) = L_{t-1}(\mathbf{x}) + \ell_t(\mathbf{x}) \ \forall \mathbf{x} \in \mathcal{C}$ .
6: end for

```

$$\Lambda(\mathbf{q}_t, \mathbf{p}, E) = \sup_{\mathbf{p}^*, f} \mathbb{E}_{\mathbf{z} \sim \mathbf{p}} \left[ \underbrace{\langle \mathbf{p} - \mathbf{p}^*, f \rangle}_{\text{true loss}} - \underbrace{\langle \mathbf{q} - \mathbf{p}^*, E(\mathbf{z}, f(\mathbf{z})) \rangle}_{\text{surrogate loss}} + \underbrace{\frac{1}{\eta} V_{\mathbf{q}_t}(\eta E(\mathbf{z}, f(\mathbf{z})))}_{\text{variance}} \right]$$

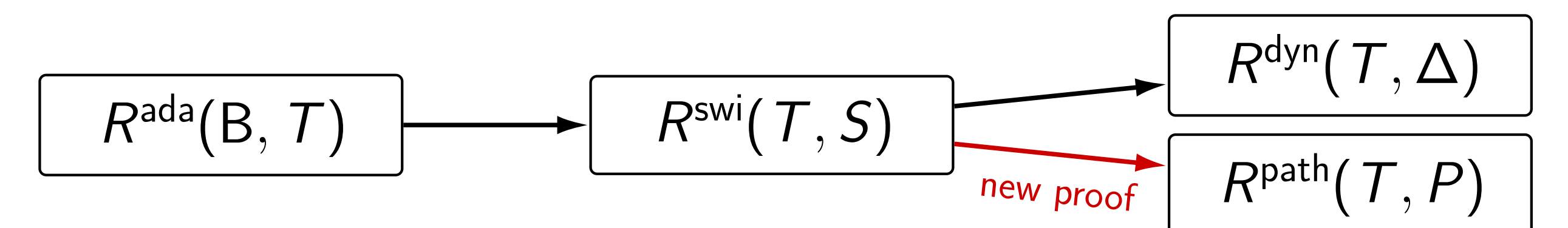
Solving  $\arg \min_{\mathbf{p}, E} \Lambda(\mathbf{q}_t, \mathbf{p}, E)$  in ExO step is statistically-sharp but exponential in computational complexity

**Theorem:** For known  $B$ ,  $R^{\text{ada}}(B, T) \lesssim d^{\frac{5}{2}}\sqrt{B}$  (GC)

## Side result: Conversions between regrets

Define **adaptive regret**: for  $B \in [T]$ ,

$$R^{\text{ada}}(B, T) := \max_{\substack{p, q \in [T], \\ 0 < q - p \leq B}} \max_{\mathbf{u} \in \Theta} \sum_{t=p}^q \mathbb{E} [f_t(\mathbf{z}_t) - f_t(\mathbf{u})].$$



Legend:  $R_1 \longrightarrow R_2$  means that if regret  $R_1$  is sublinear in  $T$  (or  $B$ ), then regret  $R_2$  is also sublinear in  $T$ , by tuning  $B$  based on  $S, \Delta, P$ .

## References

- [1] Omar Besbes, Yonatan Gur, and Assaf Zeevi. Non-stationary stochastic optimization. *Operations Research*, 63(5):1227–1244, 2015.
- [2] Yining Wang. On adaptivity in nonstationary stochastic optimization with bandit feedback. *Operations Research*, 73(2):819–828, 2025.
- [3] Peng Zhao, Guanghui Wang, Lijun Zhang, and Zhi-Hua Zhou. Bandit convex optimization in non-stationary environments. *Journal of Machine Learning Research*, 22(125):1–45, 2021.
- [4] Tianyi Chen and Georgios B Giannakis. Bandit convex optimization for scalable and dynamic IoT management. *IEEE Internet of Things Journal*, 6(1):1276–1286, 2018.
- [5] Lijun Zhang, Guanghui Wang, Wei-Wei Tu, Wei Jiang, and Zhi-Hua Zhou. Dual adaptivity: a universal algorithm for minimizing the adaptive regret of convex functions. In *International Conference on Neural Information Processing Systems*, 2021.
- [6] Guanghui Wang, Shiyin Lu, and Lijun Zhang. Adaptivity and optimality: A universal algorithm for online convex optimization. In *Proceedings of Uncertainty in Artificial Intelligence Conference*, volume 115, pages 659–668. PMLR, 2020.
- [7] Tim van Erven, Wouter M. Koolen, and Dirk van der Hoeven. Metagrad: Adaptation using multiple learning rates in online learning. *Journal of Machine Learning Research*, 22(161):1–61, 2021.
- [8] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.