

Summary and Contributions

- **Batched Algorithm for Multinomial Logistic Bandits: B-MNL-CB**
 - **Setting:** Stochastic Contextual
 - **Regret:** $\tilde{O}(\sqrt{T})$
 - **Policy Updates:** $\Omega(\log \log T)$
- **Rarely-Switching Algorithm for Multinomial Logistic Bandits: RS-MNL**
 - **Setting:** Adversarial Contextual
 - **Regret:** $\tilde{O}(\sqrt{T})$
 - **Policy Updates:** $\tilde{O}(\log T)$ (improves Sawarni et al. 2024 by logarithmic factors)
 - **Empirical Performance:** Achieves regret better than, or atleast comparable to state-of-the-art logistic and multinomial logistic bandit algorithms.

(Multinomial) Logistic Bandits

- $K + 1$ possible outcomes where the probability distribution for the outcomes is given by:

$$\mathbb{P}\{y_t = i \mid \mathbf{x}_t, \mathcal{F}_t\} = \begin{cases} z_i(\mathbf{x}_t, \boldsymbol{\theta}^*) = \frac{\exp(\mathbf{x}_t^\top \boldsymbol{\theta}_i^*)}{1 + \sum_{j=1}^K \exp(\mathbf{x}_t^\top \boldsymbol{\theta}_j^*)} & 1 \leq i \leq K, \\ z_0(\mathbf{x}_t, \boldsymbol{\theta}^*) = \frac{1}{1 + \sum_{j=1}^K \exp(\mathbf{x}_t^\top \boldsymbol{\theta}_j^*)}, & i = 0, \end{cases}$$

- Hidden Optimal Parameter: $\boldsymbol{\theta}^* = (\boldsymbol{\theta}_1^{*\top}, \dots, \boldsymbol{\theta}_K^{*\top})^\top \in \mathbb{R}^{dK}$ such that $\|\boldsymbol{\theta}^*\| \leq S$.
- Known Reward Vector: $\boldsymbol{\rho}$ such that $\|\boldsymbol{\rho}\| \leq R$ and $\rho_0 = 0$.
- Link Function $\mathbf{z}(\mathbf{x}, \boldsymbol{\theta}) = (z_1(\mathbf{x}, \boldsymbol{\theta}), \dots, z_K(\mathbf{x}, \boldsymbol{\theta}))$.
- Gradient of Link Function $\mathbf{A}(\mathbf{x}, \boldsymbol{\theta}) = \text{diag}(\mathbf{z}(\mathbf{x}, \boldsymbol{\theta})) - \mathbf{z}(\mathbf{x}, \boldsymbol{\theta})\mathbf{z}(\mathbf{x}, \boldsymbol{\theta})^\top$.
- Non-linearity parameter κ :

$$\kappa = \sup \left\{ \frac{1}{\lambda_{\min}(\mathbf{A}(\mathbf{x}, \boldsymbol{\theta}))} : \mathbf{x} \in \mathcal{X}_1 \cup \dots \cup \mathcal{X}_T, \boldsymbol{\theta} \in \Theta \right\}$$

(Distributional) Optimal Designs

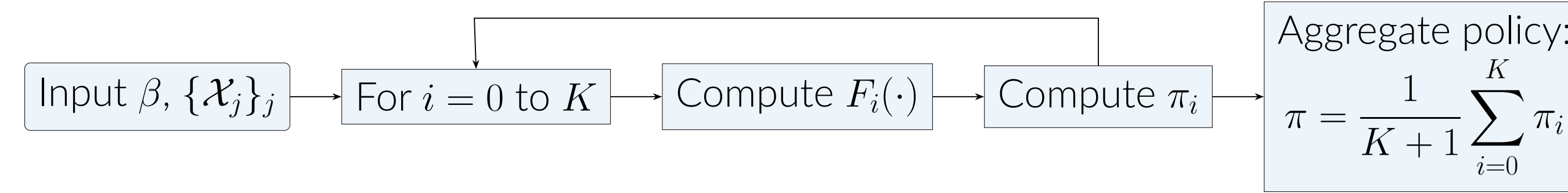
- **G-Optimal Design** π_G :
$$\max_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|_{\mathbf{V}(\pi_G)^{-1}}^2 \leq d, \quad \text{where} \quad \mathbf{V}(\pi) = \mathbb{E}_{\mathbf{x} \sim \pi}[\mathbf{x}\mathbf{x}^\top].$$
- Ruan et al. 2021 introduced **distributional optimal designs** to improve this bound
$$\mathbb{P} \left(\mathbb{E}_{\mathcal{X} \sim \mathcal{D}} \left[\max_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|_{\mathbf{V}^{-1}} \right] \leq O(\sqrt{d \log d}) \right) \geq 1 - \delta(d)$$

Batched Multinomial Contextual Bandit Algorithm: B-MNL-CB

1. **Stochastic Contextual Setting:**
 - Each arm set \mathcal{X} is sampled from the same (unknown) distribution \mathcal{D} .
 - Goal: minimize $R(T) = \mathbb{E} \left[\sum_{t=1}^T \left[\max_{\mathbf{x} \in \mathcal{X}_t} \boldsymbol{\rho}^\top \mathbf{z}(\mathbf{x}, \boldsymbol{\theta}^*) - \boldsymbol{\rho}^\top \mathbf{z}(\mathbf{x}_t, \boldsymbol{\theta}^*) \right] \right]$
2. **Algorithmic Skeleton:**

For each batch,

 - Successive Elimination: $\mathcal{X} \leftarrow \{\mathbf{x} \in \mathcal{X} : \text{UCB}(\mathbf{x}) \geq \max_{\mathbf{x} \in \mathcal{X}} \text{LCB}(\mathbf{x})\}$
 - Learner's Play: $\mathbf{x}_t \sim \pi(\mathcal{X})$
 - Policy Updates: Divide batch indices into two disjoint sets C and D :
 - $C \rightarrow$ update estimates.
 - $D \rightarrow$ parameters.
3. **Policy Calculation**



4. **Directionally Scaled Sets:** Form K different sets $F_i(\{\mathcal{X}_t\}_{t \in D}, \beta) \forall i \in [K]$.

$$F_i(\{\mathcal{X}_t\}_{t \in D}, \beta) = \left\{ \left\{ \frac{\mathbf{A}(\mathbf{x}, \hat{\boldsymbol{\theta}}_\beta)^{\frac{1}{2}}}{\sqrt{B_\beta(\mathbf{x})}} \mathbf{e}_i \otimes \mathbf{x} : \mathbf{x} \in \mathcal{X}_t \right\} : t \in D \right\}.$$

5. **Regret Guarantee:** With high probability, at the end of T rounds, the regret incurred by B-MNL-CB is bounded by

$$R(T) = \tilde{O} \left(RS^{5/4} K^{5/2} d \sqrt{T} \right)$$

Proof Sketch

- Using ideas from Zhang et al. 2023, we can decompose the regret as

$$R(T) \leq 4 \sum \mathbb{E} \left[\max_{\mathbf{x} \in \mathcal{X}} \epsilon_1(\mathbf{x}) + \max_{\mathbf{x} \in \mathcal{X}} \epsilon_2(\mathbf{x}) \right].$$

- We define $\tilde{\mathbf{X}} = \frac{\mathbf{A}(\mathbf{x}, \boldsymbol{\theta})^{1/2}}{B(\mathbf{x})} \otimes \mathbf{x}$ such that

$$\tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top = \sum_{i=1}^K \left\{ \frac{\mathbf{A}(\mathbf{x}, \hat{\boldsymbol{\theta}})^{\frac{1}{2}}}{\sqrt{B(\mathbf{x})}} \mathbf{e}_i \otimes \mathbf{x} \right\} \left\{ \frac{\mathbf{A}(\mathbf{x}, \hat{\boldsymbol{\theta}})^{\frac{1}{2}}}{\sqrt{B(\mathbf{x})}} \mathbf{e}_i \otimes \mathbf{x} \right\}^\top.$$

We use this to combine the optimal designs learned for the K different sets.

Rarely-Switching Multinomial Bandit Algorithm: RS-MNL

1. **Adversarial Contextual Setting:**
 - No assumptions on the arm sets.
 - Goal: minimize $R(T) = \sum_{t=1}^T \left[\max_{\mathbf{x} \in \mathcal{X}_t} \boldsymbol{\rho}^\top \mathbf{z}(\mathbf{x}, \boldsymbol{\theta}^*) - \boldsymbol{\rho}^\top \mathbf{z}(\mathbf{x}_t, \boldsymbol{\theta}^*) \right]$.
2. **Algorithmic Skeleton:**

At each round t

 - Check if $\det \mathbf{H}_t \geq 2 \det \mathbf{H}_\tau$ (τ : previous switching round)
 - If true, update estimate of parameters.
 - Choose arm with maximum UCB to play.
3. **Improvement over Sawarni et al. 2024:**
 - Gets rid of a warm-up switching criterion.
 - Number of switches: $\tilde{O}(\log^2 T) \leftarrow \tilde{O}(\log T)$.
4. **Regret Guarantee:** With high probability, after T rounds, the regret of RS-MNL can be bounded by:

$$R_T \leq \tilde{O} \left(RK^{3/2} S^{5/4} d \sqrt{T} \right).$$

5. **Proof Sketch:**

- Using ideas from Zhang et al. 2023, we can decompose the regret as

$$R(T) \leq \sum_{t=1}^T 2\epsilon_1(\mathbf{x}_t) + 2\epsilon_2(\mathbf{x}_t)$$

where $\epsilon_1(\mathbf{x})$ and $\epsilon_2(\mathbf{x})$ are the error terms that appear in the UCB expression.

- Bound $\epsilon_1(\mathbf{x})$ and $\epsilon_2(\mathbf{x})$ similar to B-MNL-CB (using directionally scaled sets).

Experiments

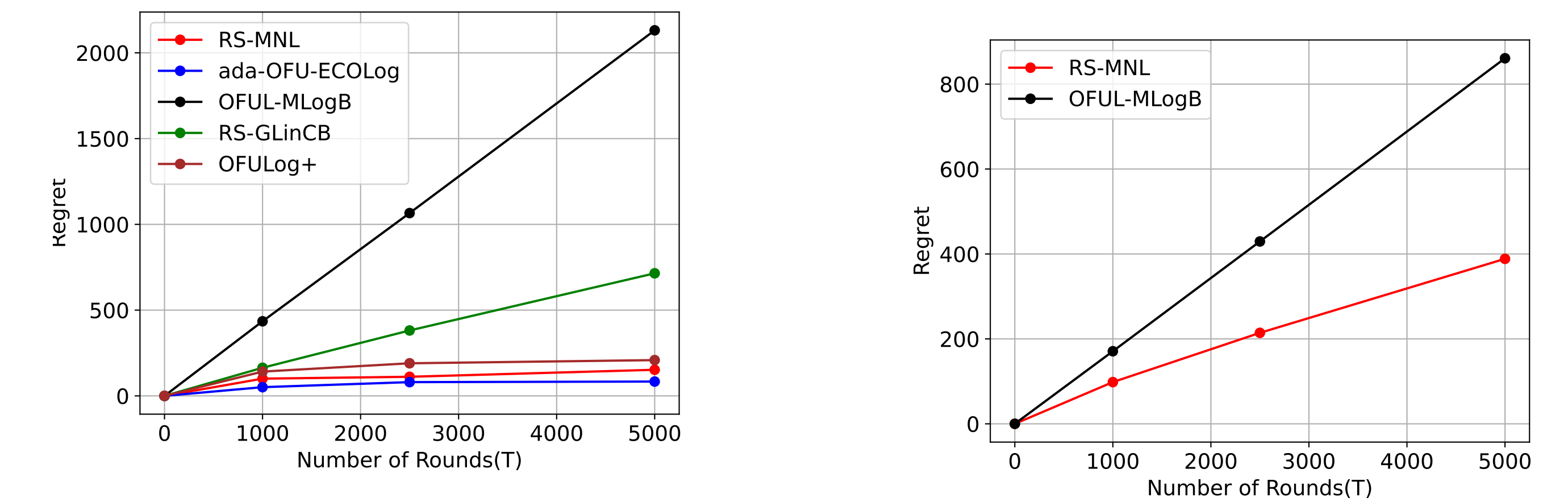


Figure 1. Regret vs. T : Logistic Setting

Regret vs. T : $K = 3$