

Improved Regret Bounds of (Multinomial) Logistic Bandits via Regret-to-Confidence-Set Conversion

Junghyun Lee¹, Se-Young Yun¹, and Kwang-Sung Jun²

¹ Kim Jaechul Graduate School of AI, KAIST, ² Department of Computer Science, University of Arizona
 {jh_lee00, yunseyoung}@kaist.ac.kr, kjun@cs.arizona.edu



Contributions

- We propose **regret-to-confidence-set conversion (R2CS)**, a new framework for converting *achievable* online learning regret bound to a confidence sequence, *without ever running the algorithm!*
- We apply R2CS to obtain the *tightest confidence set* for (multinomial) logistic losses, leading to the state-of-the-art regret guarantees for (multinomial) logistic bandits!
- Our confidence set is also numerically tight, leading to the best numerical regret by a large margin.

Logistic Bandits

Problem Setting

For $t \in [T]$:

- 1 The learner observes a potentially infinite (contextual) arm-set $\mathcal{X}_t \subset \mathbb{R}^d$
- 2 The learner chooses $\mathbf{x}_t \in \mathcal{X}_t$ according to some policy
- 3 Receive a *binary* reward $r_t | \mathbf{x}_t \sim \text{Ber}(\mu(\langle \mathbf{x}_t, \boldsymbol{\theta}_* \rangle))$,
 - $\boldsymbol{\theta}_* \in \mathbb{R}^d$ is unknown
 - $\mu(z) = (1 + e^{-z})^{-1}$ is the logistic function

Goal. Minimize:

$$\text{Reg}^B(T) := \sum_{t=1}^T \{\mu(\langle \mathbf{x}_{t,*}, \boldsymbol{\theta}_* \rangle) - \mu(\langle \mathbf{x}_t, \boldsymbol{\theta}_* \rangle)\},$$

where $\mathbf{x}_{t,*} := \arg \max_{\mathbf{x} \in \mathcal{X}_t} \mu(\langle \mathbf{x}, \boldsymbol{\theta}_* \rangle)$.

Applications. Discrete-valued rewards in interactive machine learning (e.g., clicks in news recommendations; Li et al. [2010])

Standard assumptions [Abeille et al., 2021]:

- **Assumption 1.** $\mathcal{X}_t \subseteq \mathcal{B}^d(1)$ for all $t \geq 1$.
- **Assumption 2.** $\boldsymbol{\theta}_* \in \mathcal{B}^d(S)$ with known $S > 0$.

We define the following problem-dependent quantities:

$$\kappa_*(T) := \left(\frac{1}{T} \sum_{t=1}^T \mu(\mathbf{x}_{t,*}^\top \boldsymbol{\theta}_*) \right)^{-1}, \quad \kappa_{\mathcal{X}}(T) := \max_{t \in [T]} \max_{\mathbf{x} \in \mathcal{X}_t} \frac{1}{\mu(\mathbf{x}^\top \boldsymbol{\theta}_*)},$$

and $\kappa(T) := \max_{t \in [T]} \max_{\mathbf{x} \in \mathcal{X}_t} \max_{\boldsymbol{\theta} \in \mathcal{B}^d(S)} \frac{1}{\mu(\mathbf{x}^\top \boldsymbol{\theta})}$.

These can scale *exponentially* in S !

Prior Regret Guarantees

Regret lower bound:

Theorem 2 of Abeille et al. [2021]. Let $\mathcal{X}_t = \mathcal{S}^d(1)$. Then, for any problem instance $\boldsymbol{\theta}_*$ and $T \geq d^2 \kappa_*(\boldsymbol{\theta}_*)$, there exists ϵ_T such that:

$$\min_{\pi: \text{policy}} \max_{\|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2 \leq \epsilon_T} \mathbb{E}[\text{Reg}_{\boldsymbol{\theta}, \pi}^B(T)] \geq \Omega\left(d \sqrt{\frac{T}{\kappa_*(\boldsymbol{\theta}_*)}}\right).$$

Regret upper bounds:

- **OFULog** [Abeille et al., 2021]: *Non-convex* confidence-set based UCB algorithm

$$dS^{\frac{3}{2}} \sqrt{\frac{T}{\kappa_*(T)}} + \min \{d^2 S^3 \kappa_{\mathcal{X}}(T), R_{\mathcal{X}}(T)\}.$$

- **OFULog-r** [Abeille et al., 2021]: Convex, loss-based confidence-set based UCB algorithm

$$dS^{\frac{5}{2}} \sqrt{\frac{T}{\kappa_*(T)}} + \min \{d^2 S^4 \kappa_{\mathcal{X}}(T), R_{\mathcal{X}}(T)\}.$$

- **ada-OFU-ECOLog** [Fauray et al., 2022]: Online Newton step-based algorithm

$$dS \sqrt{\frac{T}{\kappa_*(T)}} + d^2 S^6 \kappa(T).$$

Questions

- Can we construct tighter *convex, loss-based confidence set*, with improved dependency on S ?
- Can this lead to a UCB algorithm that matches or beats **ada-OFU-ECOLog**?
- Does this lead to numerically meaningful performance?

Regret-to-Confidence-Set (R2CS)

R2CS starts by directly constructing a *likelihood loss-based* confidence set centered around the norm-constrained, unregularized maximum likelihood estimator (MLE), $\hat{\boldsymbol{\theta}}_t$:

$$\hat{\boldsymbol{\theta}}_t := \arg \min_{\|\boldsymbol{\theta}\|_2 \leq S} \left\{ \mathcal{L}_t(\boldsymbol{\theta}) \triangleq \sum_{s=1}^{t-1} \ell_s(\boldsymbol{\theta}) \right\}, \quad (1)$$

where ℓ_s is the logistic loss at time s , defined as

$$\ell_s(\boldsymbol{\theta}) := -r_s \log \mu(\langle \mathbf{x}_s, \boldsymbol{\theta} \rangle) - (1 - r_s) \log(1 - \mu(\langle \mathbf{x}_s, \boldsymbol{\theta} \rangle)).$$

Theorem 1. We have $\mathbb{P}[\forall t \geq 1, \boldsymbol{\theta}_* \in \mathcal{C}_t(\delta)]$, where

$$\mathcal{C}_t(\delta) = \left\{ \boldsymbol{\theta} \in \mathcal{B}^d(S) : \mathcal{L}_t(\boldsymbol{\theta}) - \mathcal{L}_t(\hat{\boldsymbol{\theta}}_t) \leq \beta_{t-1}(\delta)^2 \right\}, \quad (2)$$

$$\beta_t(\delta) = \sqrt{10d \log \left(\frac{St}{4d} + e \right) + 2((e-2) + S) \log \frac{1}{\delta}}. \quad (3)$$

This is a strict improvement over **OFULog-r**, which has the confidence radius of $\mathcal{O}_\delta(\sqrt{dS^3 \log t})$.

Proof of R2CS for Logistic Losses

1. Decompose ℓ_s .

To use martingale concentrations, we begin by writing

$$r_s = \mu(\langle \mathbf{x}_s, \boldsymbol{\theta}_* \rangle) + \xi_s,$$

where ξ_s is a real-valued martingale difference noise.

The proof relies on the following two crucial lemmas:

Lemma 1. The following holds for any $\boldsymbol{\theta}$:

$$\ell_s(\boldsymbol{\theta}_*) = \ell_s(\boldsymbol{\theta}) + \xi_s \langle \mathbf{x}_s, \boldsymbol{\theta} - \boldsymbol{\theta}_* \rangle - \text{KL}(\mu_s(\boldsymbol{\theta}_*), \mu_s(\boldsymbol{\theta})).$$

Lemma 2. The following holds for any $\{\tilde{\boldsymbol{\theta}}_s\}$:

$$\mathcal{L}_{t+1}(\boldsymbol{\theta}_*) - \mathcal{L}_{t+1}(\hat{\boldsymbol{\theta}}_t) \leq \text{Reg}^O(t) + \zeta_1(t) - \zeta_2(t), \quad (4)$$

where

$$\text{Reg}^O(t) := \sum_{s=1}^t \{\ell_s(\tilde{\boldsymbol{\theta}}_s) - \ell_s(\hat{\boldsymbol{\theta}}_t)\},$$

$$\zeta_1(t) := \sum_{s=1}^t \xi_s \langle \mathbf{x}_s, \boldsymbol{\theta}_* - \tilde{\boldsymbol{\theta}}_s \rangle, \quad \zeta_2(t) := \sum_{s=1}^t \text{KL}(\mu_s(\boldsymbol{\theta}_*), \mu_s(\tilde{\boldsymbol{\theta}}_s)).$$

$\text{Reg}^O(t)$ is the regret incurred by the online learning algorithm of our choice up to time t , $\zeta_1(t)$ is a sum of martingale difference sequences, and $\zeta_2(t)$ is a sum of KL's.

Proof sketch. Lemma 1 follows from the first-order Taylor expansion with integral remainder and careful terms rearranging. Lemma 2 then follows immediately.

2. Use state-of-the-art online regret for $\text{Reg}^O(t)$.

Theorem 3 of Foster et al. [2018]. There is an online logistic regression algorithm with the following regret:

$$\text{Reg}^O(t) \leq 10d \log \left(\frac{St}{2d} + e \right). \quad (5)$$

We get $d \log S$ instead of dS , for free!

3. Use time-uniform Freedman to bound $\zeta_1(t)$.

Consequence of Lemma 3. For any $\eta \in [0, \frac{1}{2S}]$, the following holds w.p. at least $1 - \delta$: for all $t \geq 1$,

$$\zeta_1(t) \leq (e-2)\eta \sum_{s=1}^t \mu(\mathbf{x}_s^\top \boldsymbol{\theta}_*) \langle \mathbf{x}_s, \boldsymbol{\theta}_* - \tilde{\boldsymbol{\theta}}_s \rangle^2 + \frac{1}{\eta} \log \frac{1}{\delta}. \quad (6)$$

4. Use information-geometry to bound $\zeta_2(t)$.

Lemma 4. $\text{KL}(\mu(z_2), \mu(z_1)) = D_m(z_1, z_2)$, where D_m is the Bregman divergence generated by $m(z) = \log(1 + e^z)$.

Combined with the self-concordant analysis [Abeille et al., 2021, Lemma 8], we obtain the following:

$$-\zeta_2(t) \leq -\frac{1}{2+2S} \sum_{s=1}^t \mu(\mathbf{x}_s^\top \boldsymbol{\theta}_*) \langle \mathbf{x}_s, \boldsymbol{\theta}_* - \tilde{\boldsymbol{\theta}}_s \rangle^2. \quad (7)$$

5. Combine everything.

Set $\eta = \frac{1}{2(e-2)+2S}$, and plug Eqn. (5), (6), and (7) into Eqn. (4). \square

OFULog+

OFULog+ is of the following form:

- 1 Obtain $\hat{\boldsymbol{\theta}}_t$ (Eqn. (1)) and $\mathcal{C}_t(\delta)$ (Theorem 1)
- 2 Solve $(\mathbf{x}_t, \boldsymbol{\theta}_t) = \arg \max_{\mathbf{x} \in \mathcal{X}_t, \boldsymbol{\theta} \in \mathcal{C}_t(\delta)} \mu(\langle \mathbf{x}, \boldsymbol{\theta} \rangle)$
- 3 Play \mathbf{x}_t , then observe/receive a reward $r_t \in \{0, 1\}$.

We then have the following *state-of-the-art* regret bound:

Theorem 3. OFULog+ attains the following regret bound with probability at least $1 - \delta$:

$$\text{Reg}^B(T) \lesssim_\delta dS \sqrt{\frac{T}{\kappa_*(T)}} + \min \{d^2 S^2 \kappa_{\mathcal{X}}(T), R_{\mathcal{X}}(T)\},$$

where $R_{\mathcal{X}}(T)$ is a term relating to the arm set geometry [Abeille et al., 2021, Section 4].

Proof novelties. Time-uniform Freedman (Lemma 3) and elliptical potential count lemma [Gales et al., 2022, Lemma 7].

Experiments

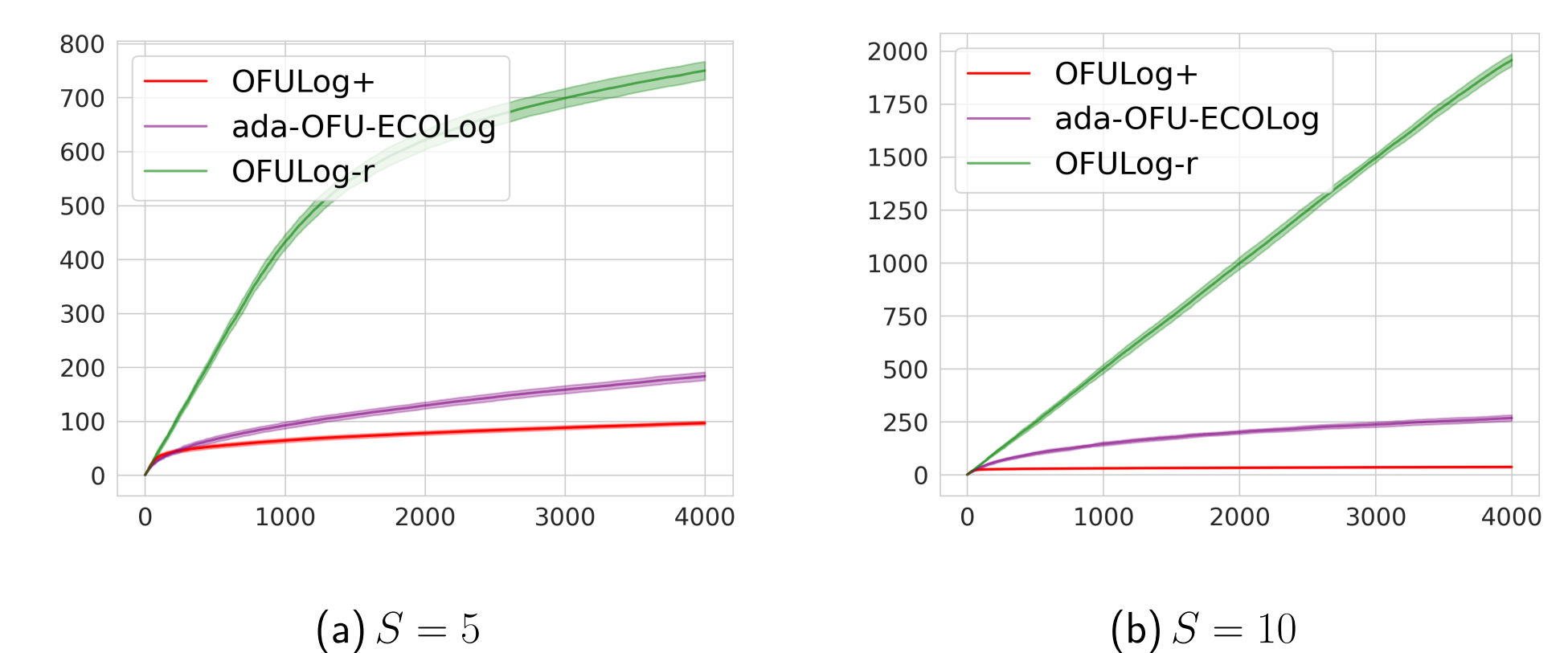


Figure 1: Numerical regrets.

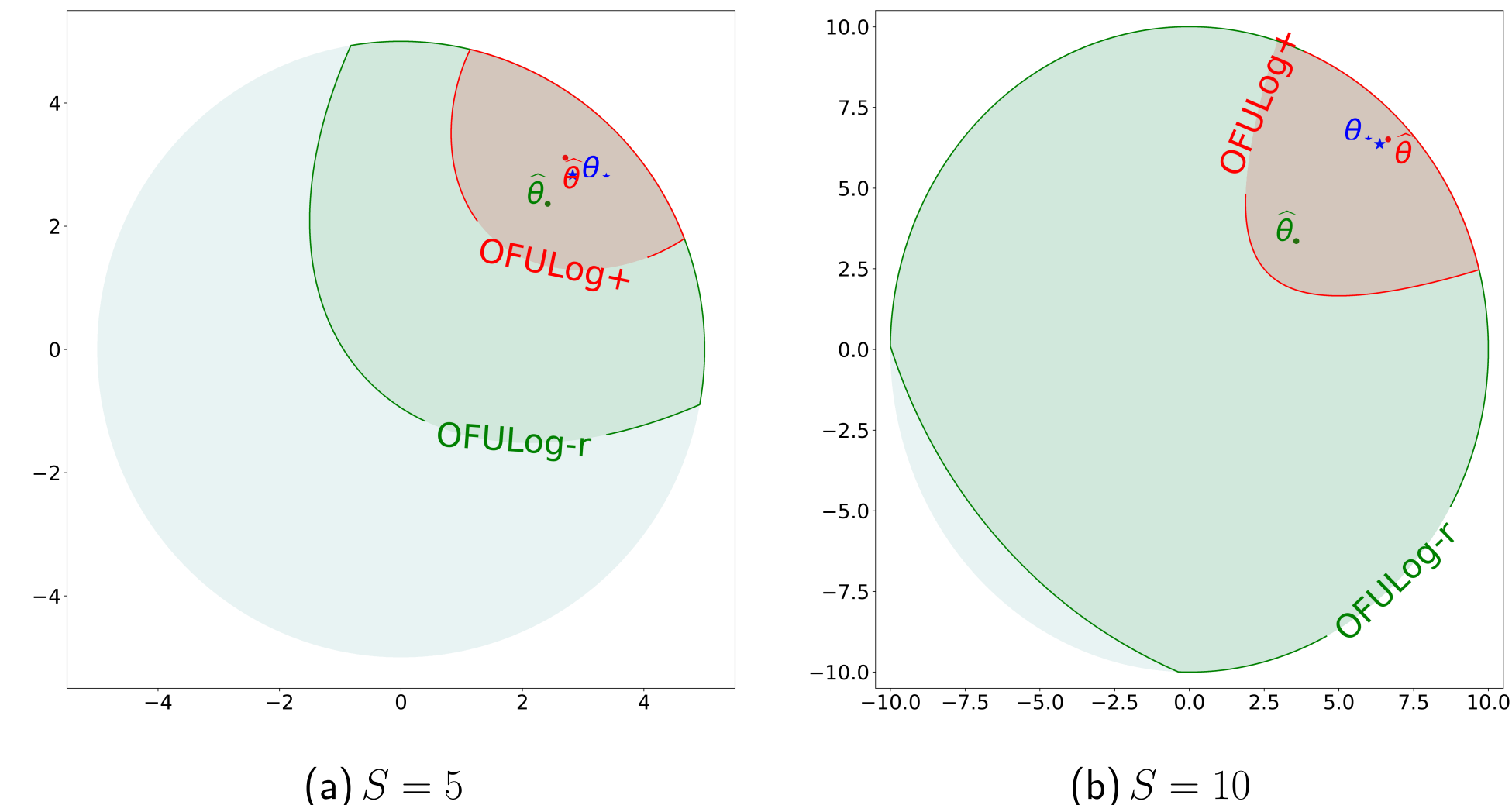


Figure 2: Confidence sets at $t = 4000$ from a single run.

Multinomial Logistic (MNL) Bandits

Via R2CS, we attain the *state-of-the-art* regret bound for MNL bandits over prior arts [Amani and Thrampoulidis, 2021, Zhang and Sugiyama, 2023]:

Theorem 5. MNL-UCB+ attains the following regret bound with probability at least $1 - \delta$:

$$\text{Reg}^B(T) \lesssim_\delta d \sqrt{KS} \min \{ \kappa(T)T, \sqrt{ST} + dKS \kappa(T) \}.$$

Open Problems

- poly(S)-free regret for (multinomial) logistic bandits?
- Extension to GLM bandits?

References

- M. Abeille et al. Instance-Wise Minimax-Optimal Algorithms for Logistic Bandits. In *AISTATS*, 2021.
- S. Amani and C. Thrampoulidis. UCB-based Algorithms for Multinomial Logistic Regression Bandits. In *NeurIPS*, 2021.
- L. Fauray et al. Jointly Efficient and Optimal Algorithms for Logistic Bandits. In *AISTATS*, 2022.
- D. J. Foster et al. Logistic Regression: The Importance of Being Improper. In *COLT*, 2018.
- S. B. Gales et al. Norm-Agnostic Linear Bandits. In *AISTATS*, 2022.
- L. Li et al. A Contextual-Bandit Approach to Personalized News Article Recommendation. In *WWW*, 2010.
- Y.-J. Zhang and M. Sugiyama. Online (Multinomial) Logistic Bandit: Improved Regret and Constant Computation Cost. In *NeurIPS*, 2023.