

## Bandit problems and Best-of-Both-Worlds Algorithms

### Combinatorial Semi-bandits

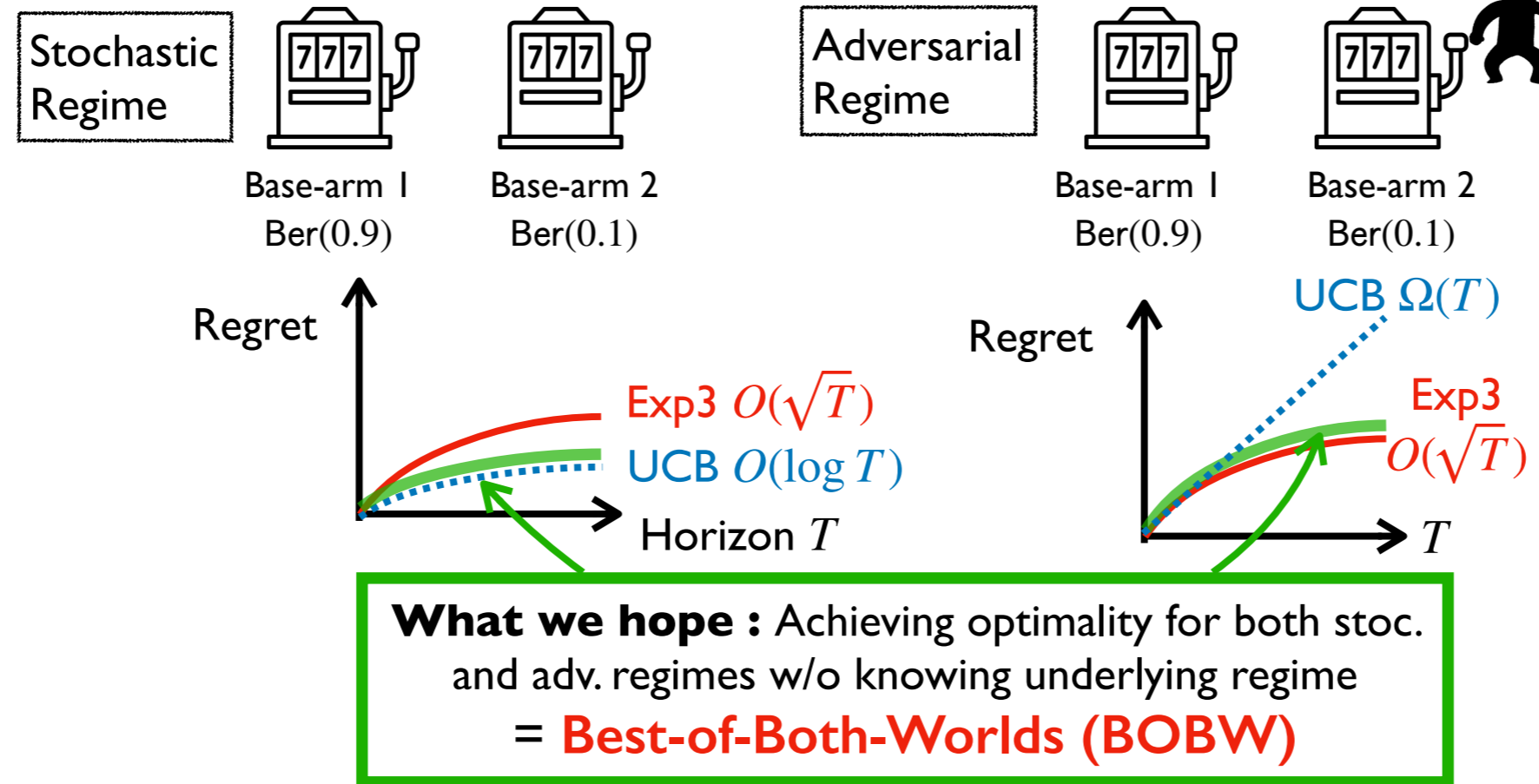
Given action set  $\mathcal{A} \subset \{0,1\}^d$   
 Adversary selects loss vectors  $\ell_1, \dots, \ell_T \in [0,1]^d$   
 At each round  $t = 1, \dots, T$ :  
 1. Learner selects  $a(t) \in \mathcal{A}$   
 2. Learner incurs a loss  $\langle \ell(t), a(t) \rangle$  and observes  $\ell_i(t)$  for  $i \in [d]$  such that  $a_i(t) = 1$

Goal: minimize regret  $R_T$  defined as

$$R_T = \mathbb{E} \left[ \sum_{t=1}^T \langle \ell(t), a(t) \rangle - \sum_{t=1}^T \langle \ell(t), a^* \rangle \right]$$

for  $a^* = \arg \min_{a \in \mathcal{A}} \sum_{t=1}^T \langle \ell(t), a \rangle$

### Stochastic and Adversarial Regime



**What we hope**: Achieving optimality for both stoc. and adv. regimes w/o knowing underlying regime = **Best-of-Both-Worlds (BOBW)**

### (Short) Research Background

There are many algorithms for BOBW algorithms mainly for multi-armed bandits [Bubeck & Slivkins 12, Zimmert & Seldin 21, etc.]

### Research Question

**Q.** Can we improve existing BOBW algorithms for semi-bandits by exploiting problem structures more?

**A.** Yes, can show regret bounds with **tight suboptimality gap** and **variance-dependent bounds!**

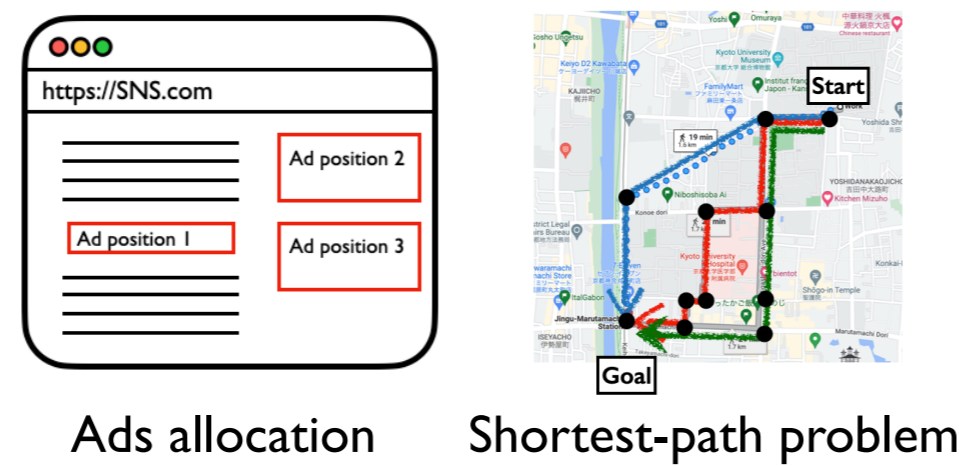
## Background and Motivation

### Why is distributional information useful?

Practically, many problems have small variances nature

- Recommender system: very small CTRs
- Shortest path problem: The required time does not vary much

Variances of each base-arm is very small  
 → Algorithms with variance-dependent bound should perform well



**Q.** Can we establish BOBW algorithms with variance-dependent bounds?

Theoretically, "Best" in BOBW literature is not the best; when  $\ell_i \sim \text{Ber}(\cdot)$

"Best" bounds in BOBW literatures

$$R_T = O\left(\sum_{i \neq i^*} \frac{\log T}{\Delta_i}\right)$$

The achievable bounds

$$R_T = \Omega\left(\sum_{i \neq i^*} \frac{\log T}{\text{KL}(\mu_i, \mu_{i^*})}\right)$$

$$\mu_i = \mathbb{E}[\ell_i]$$

$$\Delta_i = \mu_i - \mu_{i^*}$$

[Lai & Robbins 85]

**Q.** Can we fill these gaps achieving BOBW simultaneously?

### Better dependency on sub-optimality gap

Existing BOBW algorithms for combinatorial semi-bandits have a form of

$$R_T = O\left(\frac{dm \log T}{\Delta}\right)$$

for  $\Delta = \min\{\langle \mu, a - a^* \rangle : a \in \mathcal{A} \setminus \{a^*\}\}$ . [Zimmert, Luo & Wei 19, Ito 21]

**Q.** Can we achieve obtain high-resolution bounds with

arm-wise suboptimality gaps  $\Delta_{i, \min} = \min\{\langle \mu, a - a^* \rangle : a \in \mathcal{A} \setminus \{a^*\}, a_i = 1\}$ ?

## Regret Upper Bounds and Regret Analysis

### Regret Upper Bounds

$\mathcal{R}$ : bound for the stochastic regime of each algorithm

Reference	Stochastic	Adversarial	Stochastic w/ adv. corruptions
Audibert+ 14	–	$O(\sqrt{dmT})$	–
Kveton+ 15	$534 \sum_{i \in J^*} \frac{m \log T}{\Delta_{i, \min}}$	–	–
Zimmert+ 19	$O\left(\frac{dm \log T}{\Delta}\right)$	$O(\sqrt{dmT})$	$O(\mathcal{R} + \sqrt{C\mathcal{R}})$
Ito 21	$O\left(\frac{dm \log T}{\Delta}\right)$	$O(\sqrt{d \min\{L^*, Q_2, V_1\} \log T})$	$O(\mathcal{R} + \sqrt{C\mathcal{R}})$
Proposed (LS)	$\sum_{i \in J^*} \max\left\{\frac{4w\sigma_i^2}{\Delta_{i, \min}}, 2\right\} \log T$	$O(\sqrt{d \min\{L^*, Q_2\} \log T})$	$O(\mathcal{R} + \sqrt{C\mathcal{R}})$
Proposed (GD)	$\sum_{i \in J^*} \max\left\{\frac{8w\sigma_i^2}{\Delta_{i, \min}}, 4\right\} \log T$	$O(\sqrt{d \min\{L^*, Q_2, V_1\} \log T})$	$O(\mathcal{R} + \sqrt{C\mathcal{R}})$

Much better leading constant with tighter suboptimality gap  $\Delta_{i, \min}$  and variance-dependency

### Sketch of Proof

**Lemma.** In the stochastic regime w/ adversarial corruptions, it holds that

$$R_T \geq \mathbb{E} \left[ \sum_{t=1}^T \left( \frac{1}{v(\mathcal{A})} \sum_{i \in J^*} \Delta_{i, \min} (1 - a_i(t)) + \frac{1}{w(\mathcal{A})} \sum_{i \in J^*} \Delta_{i, \min} a_i(t) \right) \right] - 2Cm$$

The regret can be bounded by

$$R_T \leq O\left(\sum_{i \in J^*} \sqrt{\beta_0^2 + \frac{\sigma_i^2 P_i}{\log T}} + \sum_{i \in J^*} \sqrt{\frac{Q_i}{(\log T)^{3/2}}}\right) \text{ for } P_i = \mathbb{E}[\sum_{t=1}^T x_i(t)] \text{ and } Q_i = \mathbb{E}[\sum_{t=1}^T (1 - x_i(t))]$$

Combining the upper and lower bounds, and taking the worst case w.r.t.  $(P_i)$  and  $(Q_i)$ ,

$$\frac{R_T}{\log T} \leq 2 \frac{R_T}{\log T} - \frac{R_T}{\log T} \leq O\left(\sum_{i \in J^*} \frac{w(\mathcal{A})\sigma_i^2}{\Delta_{i, \min}} + |J^*| \frac{1}{\sqrt{\log T} \Delta_{i, \min}}\right)$$

$m = \max_{a \in \mathcal{A}} \|a\|_1$   
 $I^* = \{i \in [d] : a_i^* = 1\}$   
 $J^* = \{i \in [d] : a_i^* = 0\}$   
 $v, w \leq m$ : action-set-dep const  
 $L^* = \min_{a \in \mathcal{A}} \mathbb{E}[\sum_{t=1}^T \langle \ell(t), a \rangle]$   
 $Q_2 = \mathbb{E}[\sum_{t=1}^T \|\ell(t) - \bar{\ell}\|^2]$   
 $\bar{\ell} = T^{-1} \mathbb{E}[\sum_{t=1}^T \ell(t)]$   
 $V_1 = \mathbb{E}[\sum_{t=1}^{T-1} \|\ell(t) - \ell(t+1)\|_1]$   
 $C = \mathbb{E}[\sum_{t=1}^T \|\ell(t) - \ell'(t)\|_\infty]$

## Proposed Algorithm

### Optimistic Follow-the-Regularized-Leader (OFTRL)

- Let  $\mathcal{X}$  be a convex full of  $\mathcal{A}$
- OFTRL selects  $x_t \in \mathcal{X}$  by minimizing "predicted losses for next round + observations so far + regularizer"

"optimistic" prediction of  $\ell(t)$  + estimated losses + convex regularizer

$$x_t \in \arg \min_{x \in \mathcal{X}} \langle m(t) + \sum_{s=1}^{t-1} \hat{\ell}_s, x \rangle + \psi_t(p) \quad \hat{\ell}_s \in \mathbb{R}^d: \text{unbiased estimator of } \ell_s$$

- OFTRL with  $m(t) = 0$  is FTRL
- Most BOBW algorithms rely on (O)FTRL
- Choose  $a_t$  so that  $\mathbb{E}[a_t | x_t] = x_t$

Algorithms and analysis are similar to [Ito, Tsuchiya & Honda 22]

**Proposed Algorithms** We use OFTRL with following parameters:

**Optimistic prediction of  $\ell(t)$**

Method 1. least square (LS) estimation  $m_i(t) = \frac{1}{1 + N_i(t-1)} \left( \frac{1}{2} + \sum_{s=1}^{t-1} a_i(s) \ell_i(s) \right)$   $m(t)$  converge a.s. to  $\mu$

Method 2. gradient descent (GD) estimation → smaller leading constant

$$m_i(1) = \frac{1}{2} \text{ and } m_i(t+1) = \begin{cases} (1-\eta)m_i(t) + \eta\ell_i(t) & \text{if } i \in I(t) \\ m_i(t) & \text{otherwise} \end{cases}$$

**Loss estimation: reduced-variance estimator**

$$\hat{\ell}_i(t) = m_i(t) + \frac{a_i(t)}{x_i(t)} (\ell_i(t) - m_i(t)) \quad a_i(t): i\text{-the element of } a(t)$$

**Regularizer: base-arm-wise log-barrier + Complement of Shannon entropy**

$$\psi_t(x) = \sum_{i=1}^d \beta_i(t) \phi(x_i) \text{ with } \phi(z) = z - 1 - \log z + \log T \cdot (z + (1-z)\log(1-z))$$

with learning rate defined by

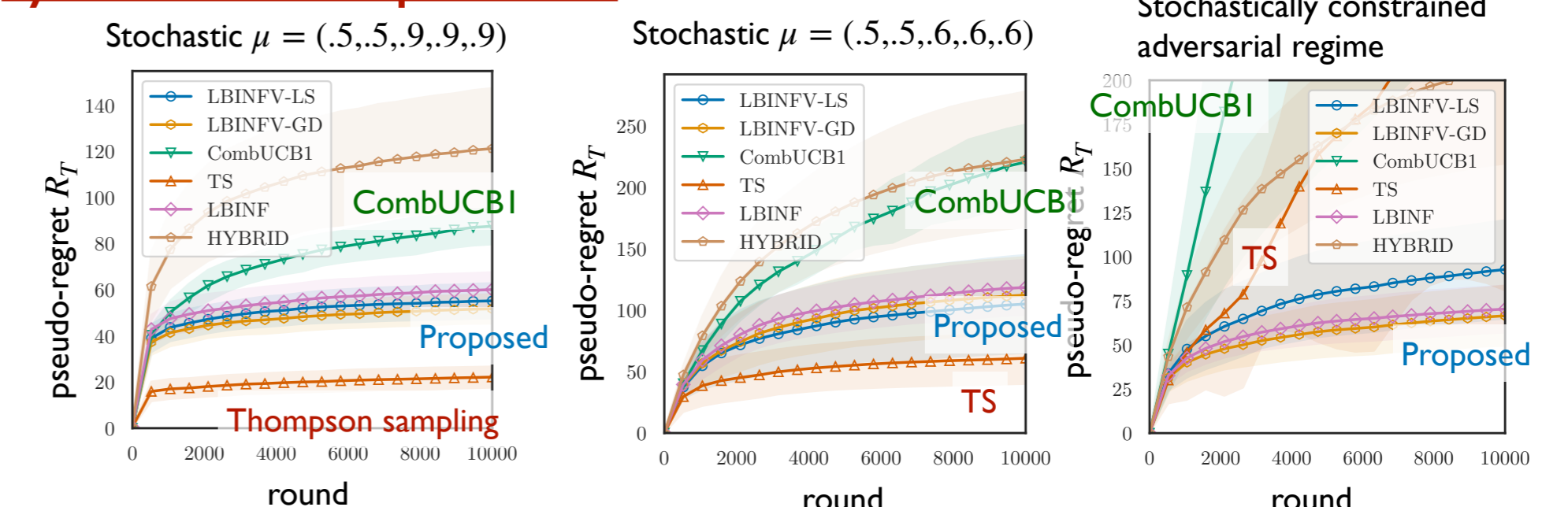
$$\beta_i(t) = \sqrt{(1+\epsilon)^2 + \frac{1}{\log T} \sum_{s=1}^{t-1} \alpha_i(s)}, \alpha_i(t) = a_i(t) (\ell_i(t) - m_i(t))^2 \min\left\{1, \frac{2(1-x_i(t))}{x_i(t)^2 \log T}\right\}$$

→  $\sigma_i^2$  as  $t \rightarrow \infty$   
 → variance-dependent bound

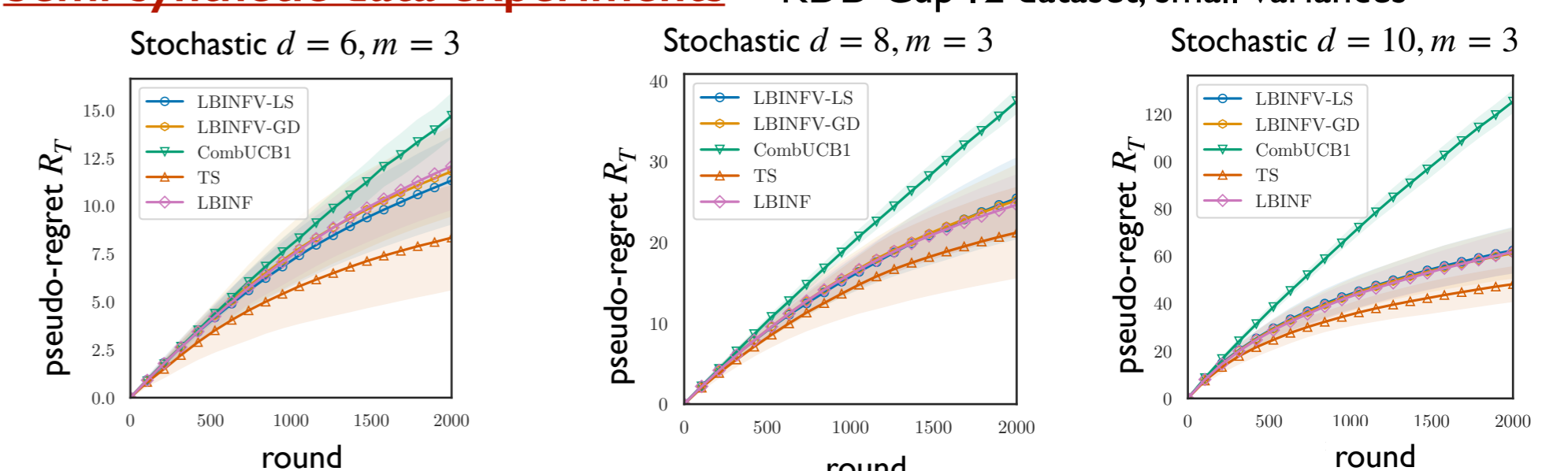
## Experiments

- $m$ -set semi-bandits with with Bernoulli base-arms
- Compare proposed algorithms (LBINFV-LS and LBINFV-GD) with CombUCB1 [Kveton, Wen, Ashkan & Szepesvári 15], Thompson sampling [Wang & Chen 18], HYBRID [Zimmert, Luo & Wei18], LBINF [Ito 21]

### Synthetic Data Experiments



### Semi-synthetic data experiments KDD Cup'12 dataset, small variances



**Proposed algorithms**

- perform better than that of existing BOBW algorithms in small-variance stoc. regime
- perform best in slightly adversarial environments where Thompson Sampling fails

S. Bubeck & A. Slivkins. "The best of both worlds: Stochastic and adversarial bandits." COLT 2012.  
 J. Zimmert & Y. Seldin. "Tallies-INF: An optimal algorithm for stochastic and adversarial bandits." JMLR 2021.  
 S. Ito, T. Tsuchiya, & J. Honda. "Adversarially robust multi-armed bandit algorithm with variance-dependent regret bounds." COLT 2022.  
 J.Y. Audibert, S. Bubeck, & G. Lugosi. "Regret in online combinatorial optimization." Mathematics of Operations Research 2014.  
 J. Kveton, Z. Wen, A. Ashkan, & Cs. Szepesvári. "Tight Regret Bounds for Stochastic Combinatorial Semi-Bandits." AISTATS 2015.  
 J. Zimmert, H. Luo, & C.Y. Wei. "Beating stochastic and adversarial semi-bandits optimally and simultaneously." ICML 2019.  
 S. Ito. "Hybrid regret bounds for combinatorial semi-bandits and adversarial linear bandits." NeurIPS 2021.  
 S. Wang & W. Chen. "Thompson sampling for combinatorial semi-bandits." ICML 2018.