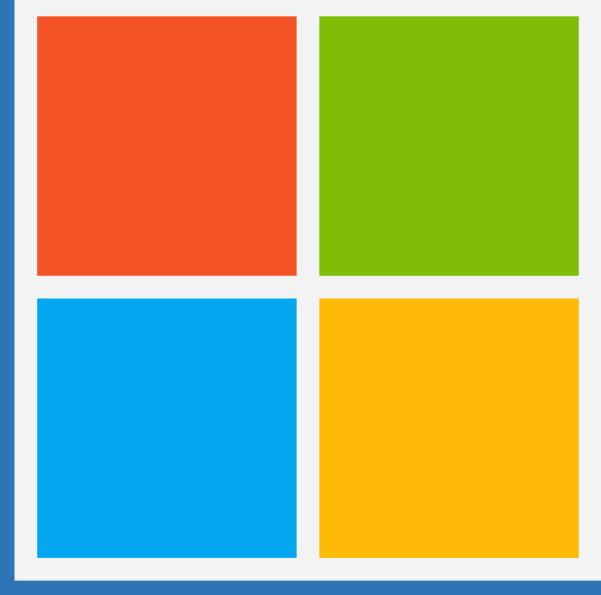


# Stochastic Online Learning with Probabilistic Graph Feedback

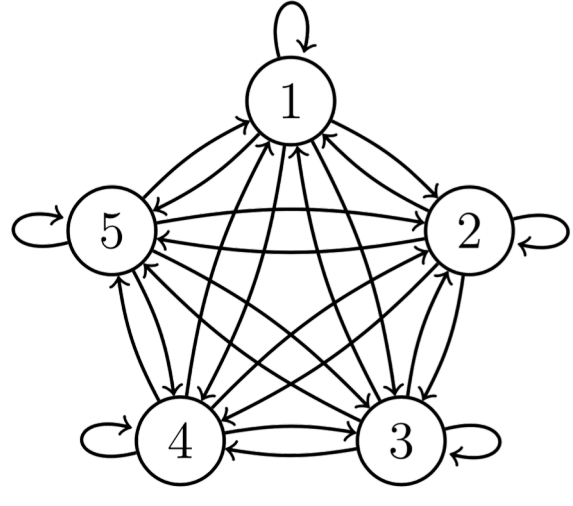


Shuai Li<sup>1</sup>; Wei Chen<sup>2</sup>; Zheng Wen<sup>3</sup>; Kwong-Sak Leung<sup>4</sup>  
<sup>1</sup>Shanghai Jiao Tong University; <sup>2</sup>Microsoft Research; <sup>3</sup>DeepMind;  
<sup>4</sup>The Chinese University of Hong Kong

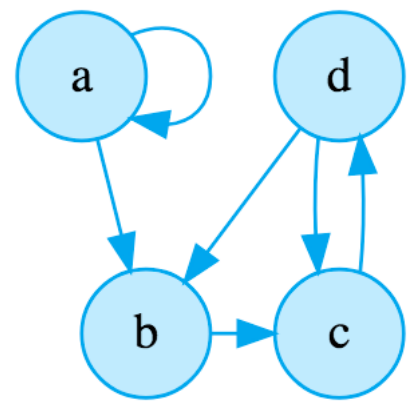


## Motivation

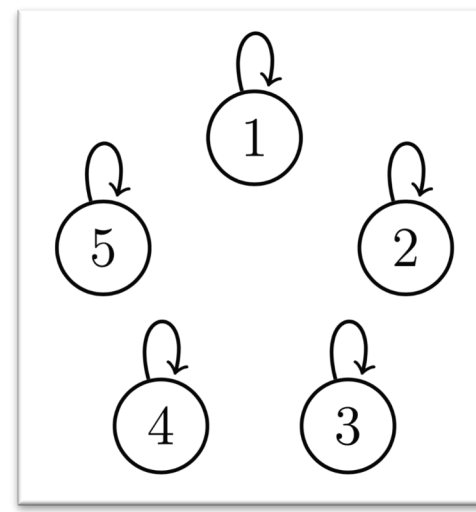
• Full information (complete graph)



• Bandit feedback (only self-loops)



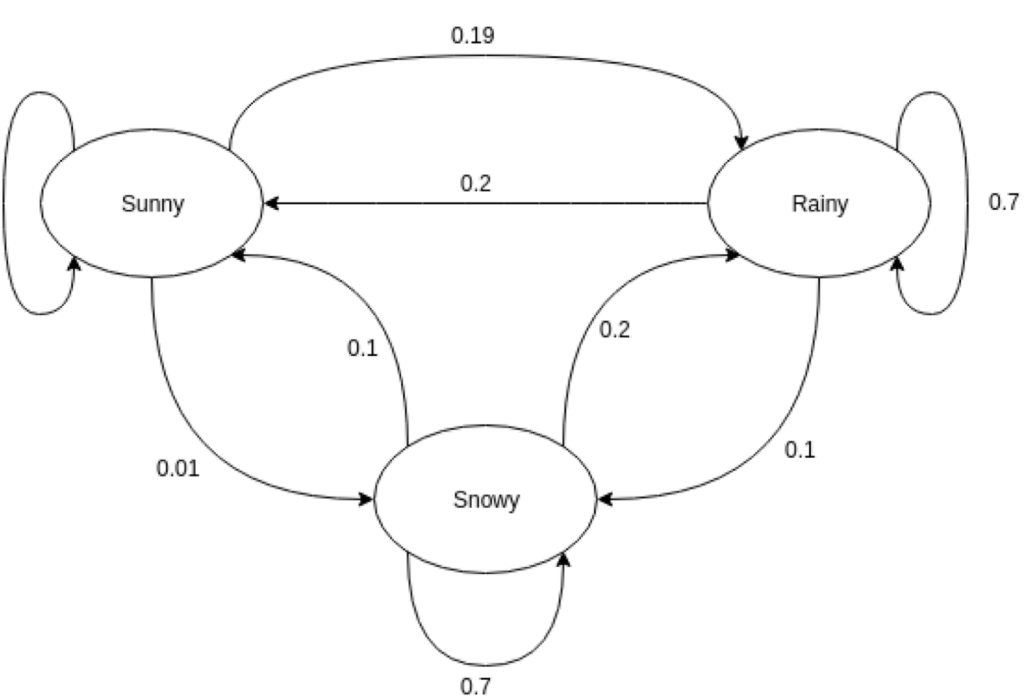
• Graph feedback can cover full information and bandit settings



• Motivating examples:

- Recommend island A would infer user's preference of island B
- Influence spread of seed A would infer influence ability of neighbor B

• Probabilistic graph



## Setting

- $V = [K]$  set of actions
- $E \subseteq V \times V$  set of directed edges
- $p: E \rightarrow (0, 1]$  triggering probabilities
- $\mu = \{\mu_i\}_{i \in V}$  unknown reward distributions
- $\theta = \{\theta_i\}_{i \in V}$  unknown reward means

• At time  $t$ ,

- The environment draws a random reward vector

$$r_t = (r_t(i): i \in V)$$

with  $r_t(i) \sim \mu_i$  and a random graph  $G_t = (V, E_t)$  with  $(i, j) \in E$  is active with probability  $p_{ij}$

- The learner selects action  $i_t \in V$  and observes

**One-Step Triggering**  
 $(j, r_t(j))$  if and only if  $(i, j) \in E_t$

**Cascade Triggering**  
 $(j, r_t(j))$  if and only if there is a path  $i \rightarrow j$  in  $G_t$

- The learner receives reward  $r_t(i_t)$

• Assumptions

- Observability:  
For each action  $j$ , there is an edge  $(i, j) \in E$  for some  $i$
- Reward distributions are of same type:  
 $KL(\mu_i, \mu_j)$  is well-defined
- Continuity:  
For each  $\mu_i, \mu_j$  and  $\epsilon$ , there exists  $\mu'_i$  such that  
 $\theta(\mu_i) + \epsilon \leq \theta(\mu'_i) \leq \theta(\mu_i) + 2\epsilon$   
 $|KL(\mu_j, \mu'_i) - KL(\mu_j, \mu_i)| \leq B\epsilon$

• Minimize the expected regret of  $T$  rounds

$$R_\mu(T; G) = T \max_{i \in V} \theta_i - \mathbb{E} \left[ \sum_{t=1}^T \theta_{i_t} \right]$$

## Lower Bounds

- Let  $i = 1$  be the best arm
- An algorithm is consistent if  $R_\mu(T) = o(T^a)$  for any  $a > 0$
- Let  $p'_{ij}$  be the probability that there is a directed path from  $i$  to  $j$  in a random realization of  $G$

**One-Step Triggering**

$$C(\mu) = \left\{ c \in [0, \infty)^V: \sum_{i \in V^{in}(1)} p_{i1} c_i \geq \frac{1}{KL(\mu_2, \mu_1)} \right. \\ \left. \sum_{i \in V^{in}(j)} p_{ij} c_i \geq \frac{1}{KL(\mu_j, \mu_1)}, \forall j \neq 1 \right\}$$

**Theorem.** For any consistent algorithm, the regret satisfies

$$\lim_{T \rightarrow \infty} \frac{R_\mu(T)}{\log T} \geq \inf_{c \in C(\mu)} \langle c, \Delta(\mu) \rangle$$

- Recovers existing works for  $p \equiv 1$

**Cascade Triggering**

$$C'(\mu) = \left\{ c \in [0, \infty)^V: \sum_{i \in V^{in}(1)} p'_{i1} c_i \geq \frac{1}{KL(\mu_2, \mu_1)} \right. \\ \left. \sum_{i \in V^{in}(j)} p'_{ij} c_i \geq \frac{1}{KL(\mu_j, \mu_1)}, \forall j \neq 1 \right\}$$

**Theorem.** For any consistent algorithm, the regret satisfies

$$\lim_{T \rightarrow \infty} \frac{R_\mu(T)}{\log T} \geq \inf_{c \in C'(\mu)} \langle c, \Delta(\mu) \rangle$$

## Results

**One-Step Triggering**

- $V^e$ : the set of exploration nodes that have the largest live probabilities among all incoming edges to some  $j$
- $p_i^e$ : the minimum exploration probability of  $i$

**Theorem.** The regret satisfies for any  $\epsilon > 0$ ,

$$R(T) = O \left( \log(T) \sum_{i=1}^K c_i(\theta, \epsilon) \Delta_i(\theta) + \log(T) \sum_{i \in V^e} \frac{\Delta_i(\theta)}{p_i^e} \right)$$

and

$$\overline{\lim}_{T \rightarrow \infty} \frac{R(T)}{\log(T)} \leq 4 \cdot \inf_{c \in C(\mu)} \langle c, \Delta(\theta) \rangle$$

holds with probability 1.

- Matched lower bound

**Cascade Triggering**

- $\hat{V}^e = \{i: p'_{ij} \geq \frac{1}{2} \max_{i'} p'_{i'j} \text{ for some } j\}$   
a relaxed version of  $V^e$
- $\hat{p}_i^e = \min \{p'_{ij}: p'_{ij} \geq \frac{1}{2} \max_{i'} p'_{i'j} \text{ for some } j\}$   
a relaxed version of  $p_i^e$

**Theorem.** The regret satisfies for any  $\epsilon > 0$ ,

$$R(T) = O \left( \sum_{i=1}^K \Delta_i(\theta) \max_{t \leq T} \{c_i(\theta, \epsilon, \eta(t)) \log t\} + \log(T) \sum_{i \in \hat{V}^e} \frac{\Delta_i(\theta)}{\hat{p}_i^e} \right)$$

and

$$\overline{\lim}_{T \rightarrow \infty} \frac{R(T)}{\log(T)} \leq 4 \cdot \inf_{c \in C'(\mu)} \langle c, \Delta(\theta) \rangle$$

holds with probability 1.

- Matched lower bound

## Algorithms

- $N_i(t)$ : selected times of action  $i$
- $M_j(t) = \sum_{i \in V^{in}(j)} N_i(t) p_{ij}$
- $n_{ij}(t)$ : observation times of action  $j$  when selecting action  $i$
- $m_j(t) = \sum_i n_{ij}(t)$

- $N^e(t)$ : number of exploration rounds on  $\theta$

**One-Step Triggering**

- $N^e = 0; \hat{\theta} = (1, 1, \dots, 1)$
- For all  $t = 1, 2, \dots, T$  do

**Gap between probabilistic graph and realizations**

- If  $m_j < M_j/2$ , play  $i_t \in \arg\max_{i \in V^{in}(j)} p_{ij}$

**Exploitation**

- Else if  $\frac{N(t)}{\log t} \in C(\hat{\theta})$ , play  $i_t = i_1(\hat{\theta})$

**Forced exploration**

**Sublinear auxiliary function**

- Else if  $M_j < 2\beta(N^e)/K$ , play  $i_t \in V^{in}(j)$  and increase  $N^e$  by 1

**Exploration by linear programming**

- Play  $i_t = i$  such that  
 $N_i < 16c(\hat{\theta}) \log(t)$   
and increase  $N^e$  by 1

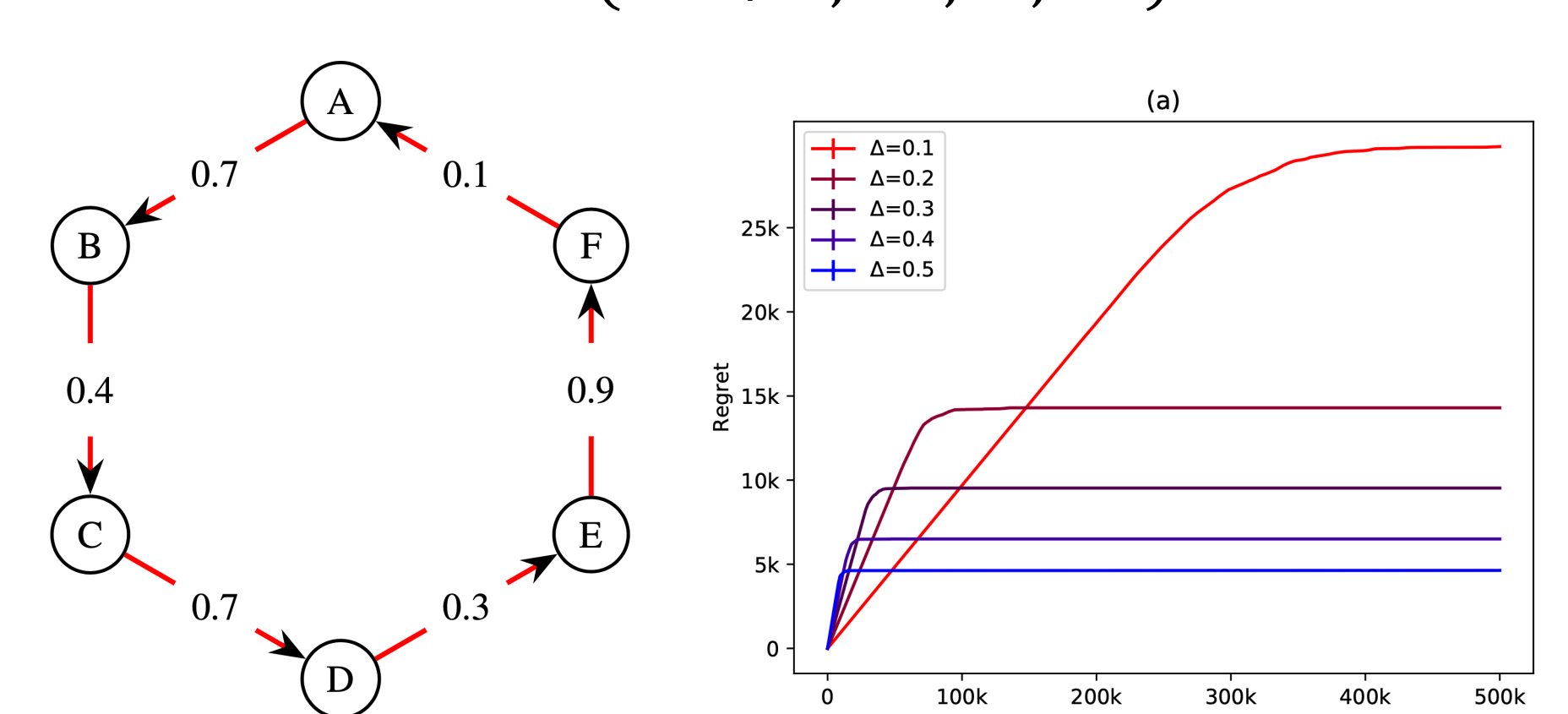
**Cascade Triggering**

- The computation of  $p'_{ij}$  is #P-hard
- Approximate them by  $p'_{ij}(t)$
- Run above algorithm on  $\tilde{G}_t$  with  $p'_{ij}(t)$
- Details omitted

## Experiments

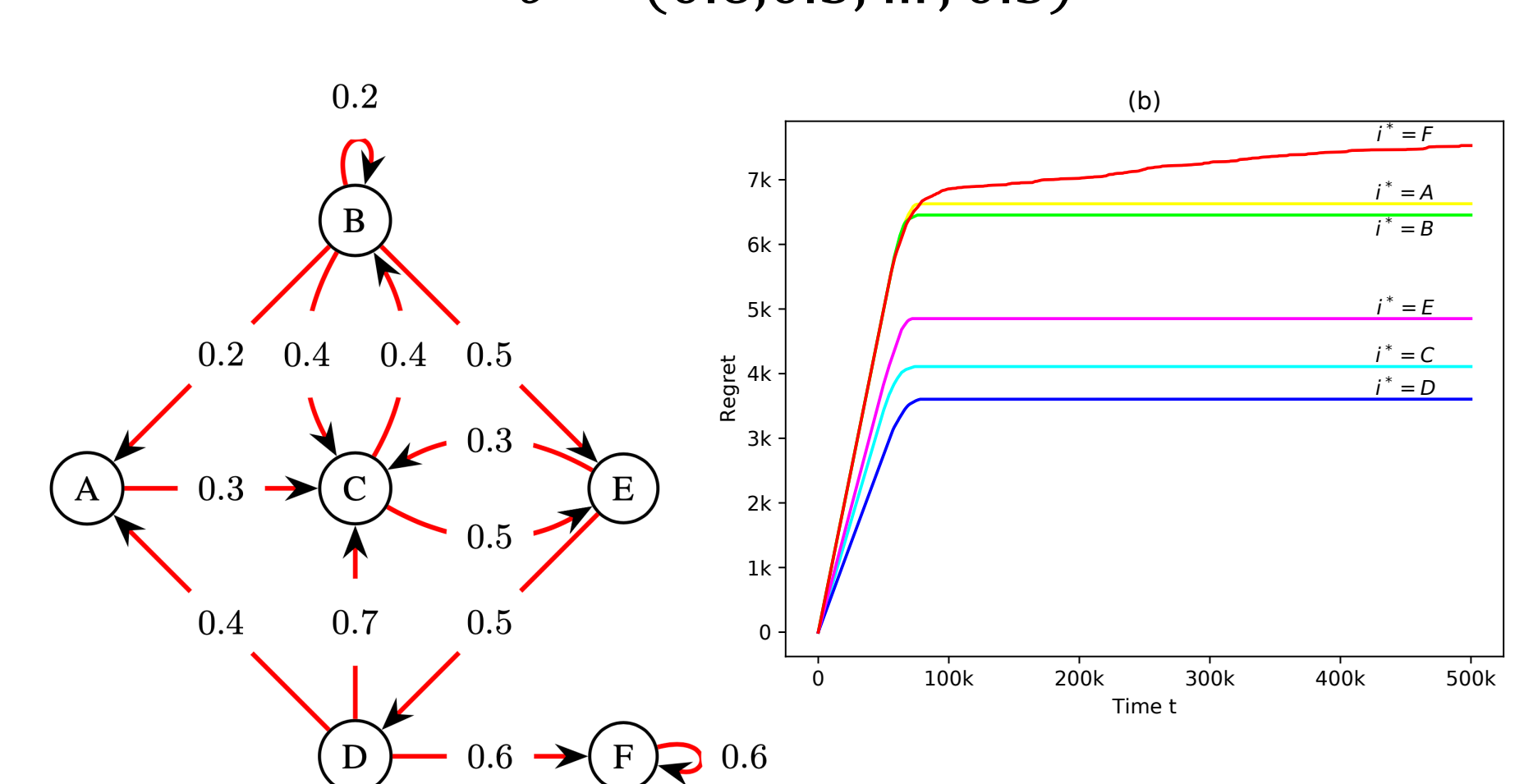
a. Cyclic feedback graph

$$\theta = (0.5 + \Delta, 0.5, \dots, 0.5)$$



b. A random feedback graph

$$\theta = (0.6, 0.5, \dots, 0.5)$$



## Conclusions

- First work, generalizes existing settings
- Also considers cascade triggering case
- Lower bound
- Upper bound, match with probability 1

## Contact

Shuai Li shuaili8@sjtu.edu.cn  
 Wei Chen weic@microsoft.com  
 Zheng Wen zhengwen@google.com  
 Kwong-Sak Leung ksleung@cse.cuhk.edu.hk

## Full Paper



## Homepage



## References

- Alon, N.; Cesa-Bianchi, N.; Dekel, O.; and Koren, T. 2015a. Online learning with feedback graphs: Beyond bandits. In Conference on Learning Theory, 23–35.
- Kocak, T.; Neu, G.; Valko, M.; and Munos, R. 2014. Efficient learning by implicit exploration in bandit problems with side observations. In Advances in Neural Information Processing Systems (NeurIPS), 613–621.
- Kocak, T.; Neu, G.; and Valko, M. 2016a. Online learning with side observations. In Uncertainty in Artificial Intelligence (UAI).
- Liu, F.; Buccapatnam, S.; and Shroff, N. 2018. Information directed sampling for stochastic bandits with graph feedback. In Thirty-Second AAAI Conference on Artificial Intelligence (AAAI).
- Mannor, S., and Shamir, O. 2011. From bandits to experts: On the value of side-observations. In Advances in Neural Information Processing Systems (NeurIPS), 684–692.
- Wu, Y.; Györfgy, A.; and Szepesvári, C. 2015. Online learning with gaussian payoffs and side observations. In Advances in Neural Information Processing Systems (NeurIPS), 1360–1368.