

# Pure Exploration with Feedback Graphs



Alessio Russo, Yichen Song and Aldo Pacchiano

---



AISTATS 2025

Pacchiano's Lab for Adaptive and Intelligent Algorithms (PLAIA)

# Best Arm Identification with Fixed Confidence

$K$  arms with reward distributions  $\nu_a$  with  $a \in \{1, \dots, K\}$ . Assume  $(\nu_a)_a$  belong to the family of **single-parameter exponential distributions**, with  $\mu_a = \mathbb{E}_{r \sim \nu_a}[r]$ .



$\mu_1$



$\mu_2$



$\mu_3$



$\mu_4$



$\mu_5$

- ▶ **Sequential:** In round  $t$  the learner pulls arm  $a_t \in [K]$  and receives the reward  $r_t \sim \nu_{a_t}$ .
- ▶ **Best Arm Identification objective:** quickly find the optimal arm  $a^* = \arg \max_a \mu_a$  with confidence  $\delta \in (0, 1/2) \Rightarrow$  **minimize sample complexity**  $\mathbb{E}[\tau]$  subject to  $\mathbb{P}(\hat{a}_\tau \neq a^*) \leq \delta$ .
  - ▶  $\tau$  is a random stopping time and  $\hat{a}_\tau$  is the estimated best arm at  $\tau$ .

Let's introduce a graph structure

# Best Arm Identification with Fixed Confidence

$K$  arms with reward distributions  $\nu_a$  with  $a \in \{1, \dots, K\}$ . Assume  $(\nu_a)_a$  belong to the family of **single-parameter exponential distributions**, with  $\mu_a = \mathbb{E}_{r \sim \nu_a}[r]$ .



$\mu_1$



$\mu_2$



$\mu_3$



$\mu_4$



$\mu_5$

- **Sequential:** In round  $t$  the learner pulls arm  $a_t \in [K]$  and receives the reward  $r_t \sim \nu_{a_t}$ .
- **Best Arm Identification objective:** quickly find the optimal arm  $a^* = \arg \max_a \mu_a$  with confidence  $\delta \in (0, 1/2) \Rightarrow$  **minimize sample complexity**  $\mathbb{E}[\tau]$  subject to  $\mathbb{P}(\hat{a}_\tau \neq a^*) \leq \delta$ .
  - $\tau$  is a random stopping time and  $\hat{a}_\tau$  is the estimated best arm at  $\tau$ .

Let's introduce a graph structure

# Best Arm Identification with Fixed Confidence

$K$  arms with reward distributions  $\nu_a$  with  $a \in \{1, \dots, K\}$ . Assume  $(\nu_a)_a$  belong to the family of **single-parameter exponential distributions**, with  $\mu_a = \mathbb{E}_{r \sim \nu_a}[r]$ .



$\mu_1$



$\mu_2$



$\mu_3$



$\mu_4$



$\mu_5$

- **Sequential:** In round  $t$  the learner pulls arm  $a_t \in [K]$  and receives the reward  $r_t \sim \nu_{a_t}$ .
- **Best Arm Identification objective:** quickly find the optimal arm  $a^* = \arg \max_a \mu_a$  with confidence  $\delta \in (0, 1/2) \Rightarrow$  **minimize sample complexity**  $\mathbb{E}[\tau]$  subject to  $\mathbb{P}(\hat{a}_\tau \neq a^*) \leq \delta$ .
  - $\tau$  is a random stopping time and  $\hat{a}_\tau$  is the estimated best arm at  $\tau$ .

Let's introduce a graph structure

# Best Arm Identification with Fixed Confidence

$K$  arms with reward distributions  $\nu_a$  with  $a \in \{1, \dots, K\}$ . Assume  $(\nu_a)_a$  belong to the family of **single-parameter exponential distributions**, with  $\mu_a = \mathbb{E}_{r \sim \nu_a}[r]$ .



$\mu_1$



$\mu_2$



$\mu_3$



$\mu_4$

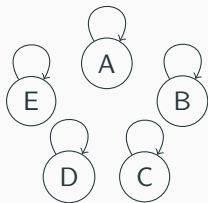


$\mu_5$

- **Sequential:** In round  $t$  the learner pulls arm  $a_t \in [K]$  and receives the reward  $r_t \sim \nu_{a_t}$ .
- **Best Arm Identification objective:** quickly find the optimal arm  $a^* = \arg \max_a \mu_a$  with confidence  $\delta \in (0, 1/2) \Rightarrow$  **minimize sample complexity**  $\mathbb{E}[\tau]$  subject to  $\mathbb{P}(\hat{a}_\tau \neq a^*) \leq \delta$ .
  - $\tau$  is a random stopping time and  $\hat{a}_\tau$  is the estimated best arm at  $\tau$ .

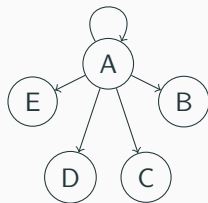
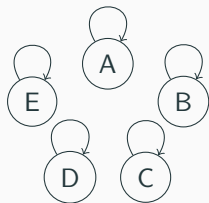
Let's introduce a graph structure

# Graph Structure



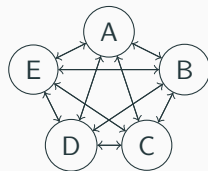
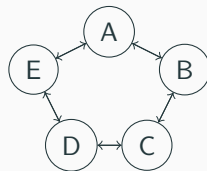
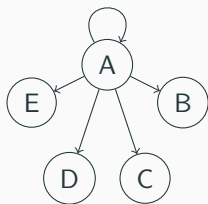
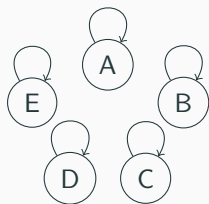
1. **Bandit model**: when selecting an action, you observe the reward of that action.
2. **Revealing action**: when selecting  $A$ , you observe the reward of all other nodes.
3. **Ring graph**: you observe only the reward of two neighboring nodes.
4. **Loopless clique**: all connected.

# Graph Structure



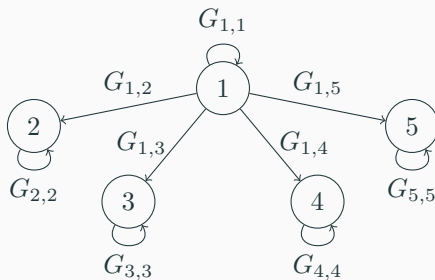
1. **Bandit model**: when selecting an action, you observe the reward of that action.
2. **Revealing action**: when selecting  $A$ , you observe the reward of all other nodes.
3. **Ring graph**: you observe only the reward of two neighboring nodes.
4. **Loopless clique**: all connected.

# Graph Structure



1. **Bandit model**: when selecting an action, you observe the reward of that action.
2. **Revealing action**: when selecting  $A$ , you observe the reward of all other nodes.
3. **Ring graph**: you observe only the reward of two neighboring nodes.
4. **Loopless clique**: all connected.



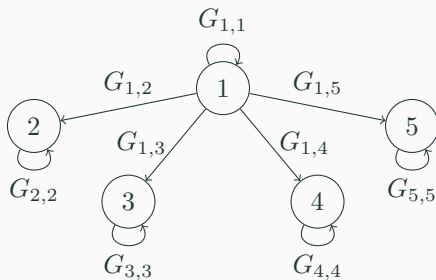


## BAI Problem

**Goal:** Estimate  $a^*$  as quickly as possible subject to  $\mathbb{P}(\hat{a}_\tau \neq a^*) \leq \delta$ .

- ▶ Graph characterized by the adjacency matrix  $G \in [0, 1]^{K \times K}$ .
- ▶ When selecting  $a$  the agent observes  $(Z_{a,u})_{u \in [K]}$ , where  $Z_{a,u} = Y_{a,u} R_u$  for all nodes  $u$ , with  $Y_{a,u} \sim \text{Ber}(G_{a,u})$  and  $R_u \sim \nu_u$ .
- ▶ What is the sample complexity lower bound? Can we use the graph to speed-up learning?

# BAI with a Graph

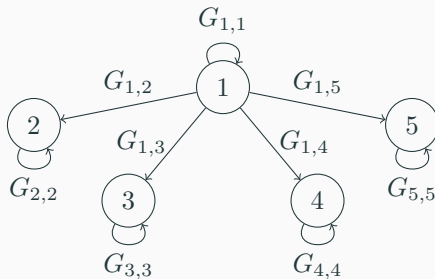


## BAI Problem

**Goal:** Estimate  $a^*$  as quickly as possible subject to  $\mathbb{P}(\hat{a}_\tau \neq a^*) \leq \delta$ .

- ▶ Graph characterized by the adjacency matrix  $G \in [0, 1]^{K \times K}$ .
- ▶ When selecting  $a$  the agent observes  $(Z_{a,u})_{u \in [K]}$ , where  $Z_{a,u} = Y_{a,u} R_u$  for all nodes  $u$ , with  $Y_{a,u} \sim \text{Ber}(G_{a,u})$  and  $R_u \sim \nu_u$ .
- ▶ What is the sample complexity **lower bound**? Can we use the graph to **speed-up** learning?

# BAI with a Graph

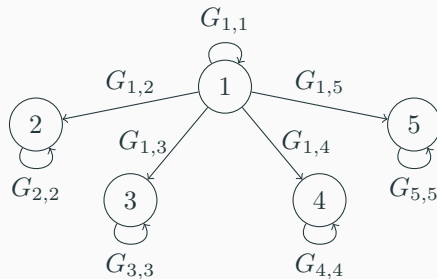


## BAI Problem

**Goal:** Estimate  $a^*$  as quickly as possible subject to  $\mathbb{P}(\hat{a}_\tau \neq a^*) \leq \delta$ .

- ▶ Graph characterized by the adjacency matrix  $G \in [0, 1]^{K \times K}$ .
- ▶ When selecting  $a$  the agent observes  $(Z_{a,u})_{u \in [K]}$ , where  $Z_{a,u} = Y_{a,u} R_u$  for all nodes  $u$ , with  $Y_{a,u} \sim \text{Ber}(G_{a,u})$  and  $R_u \sim \nu_u$ .
- ▶ What is the sample complexity **lower bound**? Can we use the graph to **speed-up** learning?

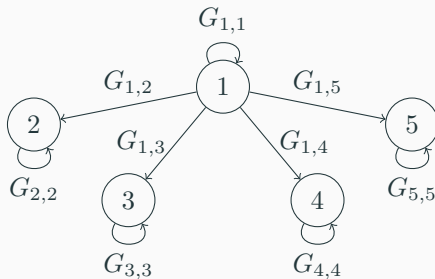
# BAI with a Graph



However, the agent **may/may not know the graph**. **Two settings:**

- ▶ **Uninformed setting:** The learner does not know  $G$  nor which edge is activated at each time-step  $t$ .
- ▶ **Informed setting:** The learner either knows  $G$  or which edge was activated after choosing a node.

# BAI with a Graph



However, the agent **may/may not know the graph**. **Two settings:**

- ▶ **Uninformed setting:** The learner does not know  $G$  nor which edge is activated at each time-step  $t$ .
- ▶ **Informed setting:** The learner either knows  $G$  or which edge was activated after choosing a node.

## Sample Complexity Lower Bounds

---

# Sample Complexity Lower Bounds - Uninformed Setting

## Theorem

For any  $\delta$ -PC algorithm and any model  $\nu$  with reward distributions  $\{\nu_u\}_{u \in V}$  with continuous support, in the **uninformed setting**<sup>1</sup> we have that

$$\mathbb{E}_\nu[\tau] \geq T^*(\nu) \log \frac{1}{2.4\delta}, \quad (1)$$

$$\underbrace{(T^*(\nu))^{-1}}_{\text{information rate}} = \underbrace{\sup_{\omega \in \Delta(V)}}_{\text{sampling policy}} \min_{u \neq a^*} (m_u + m_{a^*}) \underbrace{I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u)}_{\text{Generalized Jensen-Shannon divergence}}$$

$$\text{s.t.} \quad \underbrace{m_u}_{\text{observation rate}} = \sum_{v \in N_{in}(u)} \omega_v G_{v,u} \quad \forall u \in V.$$

Concretely, for Gaussian rewards  $\mathcal{N}(\mu_a, \lambda^2)$

$$T^*(\nu) = \inf_{\omega \in \Delta(V)} \max_{u \neq a^*} (m_u^{-1} + m_{a^*}^{-1}) \frac{2\lambda^2}{\Delta_u^2} \text{ s.t. } m = G^\top \omega \quad (\text{where } \Delta_u = \mu_{a^*} - \mu_u).$$

<sup>1</sup>**Uninformed setting:** The learner does not know  $G$  nor which edge is activated at each time-step  $t$

# Sample Complexity Lower Bounds - Uninformed Setting

## Theorem

For any  $\delta$ -PC algorithm and any model  $\nu$  with reward distributions  $\{\nu_u\}_{u \in V}$  with continuous support, in the **uninformed setting**<sup>1</sup> we have that

$$\mathbb{E}_\nu[\tau] \geq T^*(\nu) \log \frac{1}{2.4\delta}, \quad (1)$$

$$\underbrace{(T^*(\nu))^{-1}}_{\text{information rate}} = \underbrace{\sup_{\omega \in \Delta(V)}}_{\text{sampling policy}} \min_{u \neq a^*} (m_u + m_{a^*}) \underbrace{I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u)}_{\text{Generalized Jensen-Shannon divergence}}$$

$$\text{s.t.} \quad \underbrace{m_u}_{\text{observation rate}} = \sum_{v \in N_{in}(u)} \omega_v G_{v,u} \quad \forall u \in V.$$

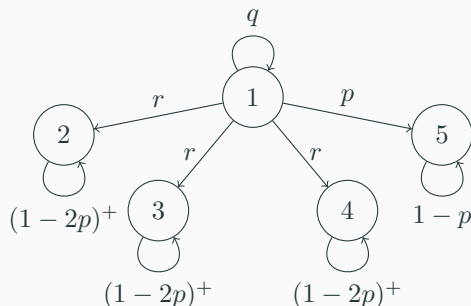
Concretely, for **Gaussian rewards**  $\mathcal{N}(\mu_a, \lambda^2)$

$$T^*(\nu) = \inf_{\omega \in \Delta(V)} \max_{u \neq a^*} (m_u^{-1} + m_{a^*}^{-1}) \frac{2\lambda^2}{\Delta_u^2} \text{ s.t. } m = G^\top \omega \quad (\text{where } \Delta_u = \mu_{a^*} - \mu_u).$$

<sup>1</sup>**Uninformed setting:** The learner does not know  $G$  nor which edge is activated at each time-step  $t$



## An Example: The Loopy Star

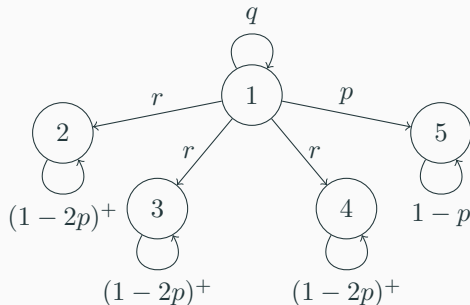


**Loopy star graph.** To each edge is associated an activation probability (obs. that  $(x)^+ = \max(x, 0)$ ).

We consider Gaussian rewards, with  $\lambda = 1$ ,  $\mu_5 = 1$  and  $\mu_u = 0.5, u \in \{1, \dots, 4\}$ .

- This graph is the union of a bandit feedback graph and revealing action graph. Removing any self-loop changes the minimax regret from  $\tilde{\Theta}(\sqrt{\alpha(G)T})$  to  $\tilde{\Theta}(T^{2/3})$  [ACBDK15].

## An Example: The Loopy Star

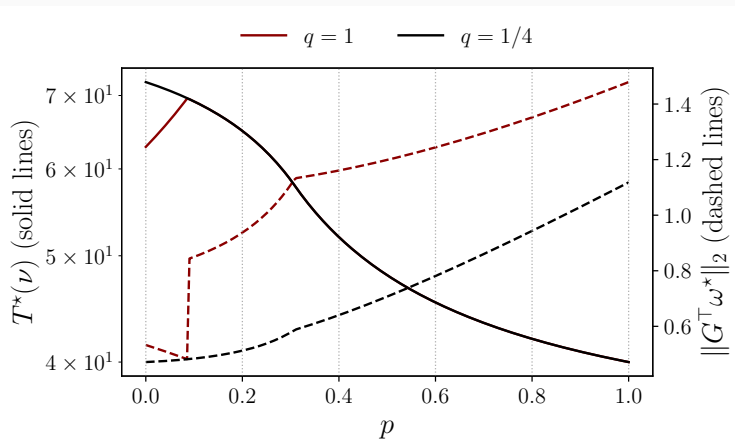


**Loopy star graph.** To each edge is associated an activation probability (obs. that  $(x)^+ = \max(x, 0)$ ).

We consider Gaussian rewards, with  $\lambda = 1$ ,  $\mu_5 = 1$  and  $\mu_u = 0.5, u \in \{1, \dots, 4\}$ .

- This graph is the union of a bandit feedback graph and revealing action graph. **Removing any self-loop changes the minimax regret from  $\tilde{\Theta}(\sqrt{\alpha(G)T})$  to  $\tilde{\Theta}(T^{2/3})$  [ACBDK15].**

# The Loopy Star



Loopy star example with  $r = 1/4$ . The solid lines depict  $T^*(\nu)$  for  $q = 1$  and  $q = 1/4$  for different values of  $p$ . Similarly, on the right axis, the dashed lines show  $\|G^\top \omega^*\|_2$ , which indicates the amount of information gathered per time-step.

# Sample Complexity Lower Bounds - Uninformed Setting [Proof 1/3]

**Overall proof idea:** take the log-likelihood ratio (LLR) of the observed data up to time  $\tau$  between the true model  $\nu$  and an alternative model  $\nu'$  that admits a different optimal vertex.

- Selecting the model  $\nu'$  that minimizes the LLR yields a **lower bound on the sample complexity**.

**Step 1 (LLR):** Consider two bandit models  $\nu = \{G, (\nu_u)_u\}$ ,  $\nu' = \{G', (\nu'_u)_u\}$ . For each  $v$ ,  $\nu_v$  and  $\nu'_v$  have, respectively, densities  $f_v$  and  $f'_v$ .  $Z_{v,u} = Y_{v,u}R_u$  has density  $f_{v,u}$  (sim.  $f'_{v,u}$ )

$$\begin{aligned} L_t &= \ln \frac{d\mathbb{P}_\nu(V_1, Z_1, \dots, V_t, Z_t)}{d\mathbb{P}_{\nu'}(V_1, Z_1, \dots, V_t, Z_t)}, && (V_t \text{ is the chosen vertex; } Z_t \text{ are the observed } \{Z_{v,u}\}_{v,u} \text{ at time } t) \\ &= \sum_{u \in V} \sum_{v \in N_{in}(u)} \sum_{j=1}^{N_v(t)} \ln \left( \frac{f_{v,u}(W_{j,(v,u)})}{f'_{v,u}(W_{j,(v,u)})} \right), && (W_{j,(v,u)} \text{ is the } j\text{-th obs. of } Z_{v,u}) \\ \implies \mathbb{E}_\nu[L_t] &= \sum_{u \in V} \sum_{v \in N_{in}(u)} \mathbb{E}_\nu[N_v(t)] \text{KL}(\nu_{v,u}, \nu'_{v,u}). \end{aligned}$$

# Sample Complexity Lower Bounds - Uninformed Setting [Proof 1/3]

**Overall proof idea:** take the log-likelihood ratio (LLR) of the observed data up to time  $\tau$  between the true model  $\nu$  and an alternative model  $\nu'$  that admits a different optimal vertex.

- Selecting the model  $\nu'$  that minimizes the LLR yields a **lower bound on the sample complexity**.

**Step 1 (LLR):** Consider two bandit models  $\nu = \{G, (\nu_u)_u\}$ ,  $\nu' = \{G', (\nu'_u)_u\}$ . For each  $v$ ,  $\nu_v$  and  $\nu'_v$  have, respectively, densities  $f_v$  and  $f'_v$ .  $Z_{v,u} = Y_{v,u}R_u$  has density  $f_{v,u}$  (sim.  $f'_{v,u}$ )

$$\begin{aligned} L_t &= \ln \frac{d\mathbb{P}_\nu(V_1, Z_1, \dots, V_t, Z_t)}{d\mathbb{P}_{\nu'}(V_1, Z_1, \dots, V_t, Z_t)}, && (V_t \text{ is the chosen vertex; } Z_t \text{ are the observed } \{Z_{v,u}\}_{v,u} \text{ at time } t) \\ &= \sum_{u \in V} \sum_{v \in N_{in}(u)} \sum_{j=1}^{N_v(t)} \ln \left( \frac{f_{v,u}(W_{j,(v,u)})}{f'_{v,u}(W_{j,(v,u)})} \right), && (W_{j,(v,u)} \text{ is the } j\text{-th obs. of } Z_{v,u}) \\ \implies \mathbb{E}_\nu[L_t] &= \sum_{u \in V} \sum_{v \in N_{in}(u)} \mathbb{E}_\nu[N_v(t)] \text{KL}(\nu_{v,u}, \nu'_{v,u}). \end{aligned}$$

# Sample Complexity Lower Bounds - Uninformed Setting [Proof 1/3]

**Overall proof idea:** take the log-likelihood ratio (LLR) of the observed data up to time  $\tau$  between the true model  $\nu$  and an alternative model  $\nu'$  that admits a different optimal vertex.

- Selecting the model  $\nu'$  that minimizes the LLR yields a **lower bound on the sample complexity**.

**Step 1 (LLR):** Consider two bandit models  $\nu = \{G, (\nu_u)_u\}$ ,  $\nu' = \{G', (\nu'_u)_u\}$ . For each  $v$ ,  $\nu_v$  and  $\nu'_v$  have, respectively, densities  $f_v$  and  $f'_v$ .  $Z_{v,u} = Y_{v,u}R_u$  has density  $f_{v,u}$  (sim.  $f'_{v,u}$ )

$$\begin{aligned} L_t &= \ln \frac{d\mathbb{P}_\nu(V_1, Z_1, \dots, V_t, Z_t)}{d\mathbb{P}_{\nu'}(V_1, Z_1, \dots, V_t, Z_t)}, && (V_t \text{ is the chosen vertex; } Z_t \text{ are the observed } \{Z_{v,u}\}_{v,u} \text{ at time } t) \\ &= \sum_{u \in V} \sum_{v \in N_{in}(u)} \sum_{j=1}^{N_v(t)} \ln \left( \frac{f_{v,u}(W_{j,(v,u)})}{f'_{v,u}(W_{j,(v,u)})} \right), && (W_{j,(v,u)} \text{ is the } j\text{-th obs. of } Z_{v,u}) \\ \implies \mathbb{E}_\nu[L_t] &= \sum_{u \in V} \sum_{v \in N_{in}(u)} \mathbb{E}_\nu[N_v(t)] \text{KL}(\nu_{v,u}, \nu'_{v,u}). \end{aligned}$$

## Sample Complexity Lower Bounds - Uninformed Setting [Proof 2/3]

Using an **information processing inequality**[KCG16], we can lower bound the expected LLR at  $\tau$  as

$$\mathbb{E}_\nu[L_\tau] \geq \log(1/(2.4\delta)),$$

and by letting  $\omega_v = \mathbb{E}_\nu[N_v(\tau)]/\mathbb{E}_\nu[\tau]$ , we obtain

$$\underbrace{\mathbb{E}_\nu[\tau]}_{\text{sample complexity}} \sum_{u \in V} \sum_{v \in N_{in}(u)} \omega_v \text{KL}(\nu_{v,u}, \nu'_{v,u}) \geq \log(1/(2.4\delta)).$$

Lastly, because of the nature of the problem, we can prove that for continuous rewards we have<sup>2</sup>

$$\text{KL}(\nu_{v,u}, \nu'_{v,u}) = \underbrace{\text{kl}(G_{v,u}, G'_{v,u})}_{\text{Bernoulli KL divergence}} + G_{v,u} \underbrace{\text{KL}(\nu_u, \nu'_u)}_{\text{KL divergence of the rewards}}.$$

---

<sup>2</sup>Recall that  $G_{v,u}$  is the edge activation probability (probability of observing  $u$  when selecting  $v$ ).

## Sample Complexity Lower Bounds - Uninformed Setting [Proof 2/3]

Using an **information processing inequality**[KCG16], we can lower bound the expected LLR at  $\tau$  as

$$\mathbb{E}_\nu[L_\tau] \geq \log(1/(2.4\delta)),$$

and by letting  $\omega_v = \mathbb{E}_\nu[N_v(\tau)]/\mathbb{E}_\nu[\tau]$ , we obtain

$$\underbrace{\mathbb{E}_\nu[\tau]}_{\text{sample complexity}} \sum_{u \in V} \sum_{v \in N_{in}(u)} \omega_v \text{KL}(\nu_{v,u}, \nu'_{v,u}) \geq \log(1/(2.4\delta)).$$

Lastly, because of the nature of the problem, we can prove that for continuous rewards we have<sup>2</sup>

$$\text{KL}(\nu_{v,u}, \nu'_{v,u}) = \underbrace{\text{kl}(G_{v,u}, G'_{v,u})}_{\text{Bernoulli KL divergence}} + G_{v,u} \underbrace{\text{KL}(\nu_u, \nu'_u)}_{\text{KL divergence of the rewards}}.$$

---

<sup>2</sup>Recall that  $G_{v,u}$  is the edge activation probability (probability of observing  $u$  when selecting  $v$ ).



## Sample Complexity Lower Bounds - Uninformed Setting [Proof 2/3]

Using an **information processing inequality**[KCG16], we can lower bound the expected LLR at  $\tau$  as

$$\mathbb{E}_\nu[L_\tau] \geq \log(1/(2.4\delta)),$$

and by letting  $\omega_v = \mathbb{E}_\nu[N_v(\tau)]/\mathbb{E}_\nu[\tau]$ , we obtain

$$\underbrace{\mathbb{E}_\nu[\tau]}_{\text{sample complexity}} \sum_{u \in V} \sum_{v \in N_{in}(u)} \omega_v \text{KL}(\nu_{v,u}, \nu'_{v,u}) \geq \log(1/(2.4\delta)).$$

Lastly, because of the nature of the problem, we can prove that for continuous rewards we have<sup>2</sup>

$$\text{KL}(\nu_{v,u}, \nu'_{v,u}) = \underbrace{\text{kl}(G_{v,u}, G'_{v,u})}_{\text{Bernoulli KL divergence}} + G_{v,u} \underbrace{\text{KL}(\nu_u, \nu'_u)}_{\text{KL divergence of the rewards}}.$$

---

<sup>2</sup>Recall that  $G_{v,u}$  is the edge activation probability (probability of observing  $u$  when selecting  $v$ ).

## Sample Complexity Lower Bounds - Uninformed Setting [Proof 3/3]

**Step 2 (Optimizing  $\nu'$ ):** We focus on **alternative models  $\nu'$**  that admit a different optimal vertex

$$\text{Alt}(\nu) = \cup_{v \neq a^*} \text{Alt}_v(\nu), \quad \text{Alt}_v(\nu) = \{\nu' \mid \mu'_v > \mu'_{a^*}\}.$$

Choose  $\nu'$  as to **minimize the LLR!**

$$\begin{aligned} & \inf_{\nu' \in \text{Alt}(\nu)} \sum_{u \in V} \sum_{v \in N_{in}(u)} \omega_v \text{KL}(\nu_{v,u}, \nu'_{v,u}) \\ &= \min_{u \neq a^*} \inf_{\nu': \mu'_u \geq \mu'_{a^*}} \sum_{v \in N_{in}(u)} \omega_v G_{v,u} \text{KL}(\nu_u, \nu'_u) + \sum_{w \in N_{in}(a^*)} \omega_w G_{w,a^*} \text{KL}(\nu_{a^*}, \nu'_{a^*}), \\ &= \min_{u \neq a^*} \inf_{\nu': \mu'_u \geq \mu'_{a^*}} m_u \text{KL}(\nu_u, \nu'_u) + m_{a^*} \text{KL}(\nu_{a^*}, \nu'_{a^*}). \quad (m_u := \sum_{v \in N_{in}(u)} \omega_v G_{v,u}) \end{aligned}$$

Therefore, by **optimizing over  $\nu'$**  as in [GK16, Lemma 3] we obtain

$$\min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \geq \log(1/(2.4\delta)).$$

## Sample Complexity Lower Bounds - Uninformed Setting [Proof 3/3]

**Step 2 (Optimizing  $\nu'$ ):** We focus on **alternative models  $\nu'$**  that admit a different optimal vertex

$$\text{Alt}(\nu) = \cup_{v \neq a^*} \text{Alt}_v(\nu), \quad \text{Alt}_v(\nu) = \{\nu' \mid \mu'_v > \mu'_{a^*}\}.$$

Choose  $\nu'$  as to **minimize the LLR!**

$$\begin{aligned} & \inf_{\nu' \in \text{Alt}(\nu)} \sum_{u \in V} \sum_{v \in N_{in}(u)} \omega_v \text{KL}(\nu_{v,u}, \nu'_{v,u}) \\ &= \min_{u \neq a^*} \inf_{\nu': \mu'_u \geq \mu'_{a^*}} \sum_{v \in N_{in}(u)} \omega_v G_{v,u} \text{KL}(\nu_u, \nu'_u) + \sum_{w \in N_{in}(a^*)} \omega_w G_{w,a^*} \text{KL}(\nu_{a^*}, \nu'_{a^*}), \\ &= \min_{u \neq a^*} \inf_{\nu': \mu'_u \geq \mu'_{a^*}} m_u \text{KL}(\nu_u, \nu'_u) + m_{a^*} \text{KL}(\nu_{a^*}, \nu'_{a^*}). \quad (m_u := \sum_{v \in N_{in}(u)} \omega_v G_{v,u}) \end{aligned}$$

Therefore, by **optimizing over  $\nu'$**  as in [GK16, Lemma 3] we obtain

$$\min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \geq \log(1/(2.4\delta)).$$

# Sample Complexity Lower Bounds - Uninformed Setting [Proof 3/3]

**Step 2 (Optimizing  $\nu'$ ):** We focus on **alternative models  $\nu'$**  that admit a different optimal vertex

$$\text{Alt}(\nu) = \cup_{v \neq a^*} \text{Alt}_v(\nu), \quad \text{Alt}_v(\nu) = \{\nu' \mid \mu'_v > \mu'_{a^*}\}.$$

Choose  $\nu'$  as to **minimize the LLR!**

$$\begin{aligned} & \inf_{\nu' \in \text{Alt}(\nu)} \sum_{u \in V} \sum_{v \in N_{in}(u)} \omega_v \text{KL}(\nu_{v,u}, \nu'_{v,u}) \\ &= \min_{u \neq a^*} \inf_{\nu': \mu'_u \geq \mu'_{a^*}} \sum_{v \in N_{in}(u)} \omega_v G_{v,u} \text{KL}(\nu_u, \nu'_u) + \sum_{w \in N_{in}(a^*)} \omega_w G_{w,a^*} \text{KL}(\nu_{a^*}, \nu'_{a^*}), \\ &= \min_{u \neq a^*} \inf_{\nu': \mu'_u \geq \mu'_{a^*}} m_u \text{KL}(\nu_u, \nu'_u) + m_{a^*} \text{KL}(\nu_{a^*}, \nu'_{a^*}). \quad (m_u := \sum_{v \in N_{in}(u)} \omega_v G_{v,u}) \end{aligned}$$

Therefore, by **optimizing over  $\nu'$**  as in [GK16, Lemma 3] we obtain

$$\min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \geq \log(1/(2.4\delta)).$$

## Sample Complexity Lower Bounds - Uninformed Setting [Proof 3/3]

**Step 2 (Optimizing  $\nu'$ ):** We focus on **alternative models  $\nu'$**  that admit a different optimal vertex

$$\text{Alt}(\nu) = \cup_{v \neq a^*} \text{Alt}_v(\nu), \quad \text{Alt}_v(\nu) = \{\nu' \mid \mu'_v > \mu'_{a^*}\}.$$

Choose  $\nu'$  as to **minimize the LLR!**

$$\begin{aligned} & \inf_{\nu' \in \text{Alt}(\nu)} \sum_{u \in V} \sum_{v \in N_{in}(u)} \omega_v \text{KL}(\nu_{v,u}, \nu'_{v,u}) \\ &= \min_{u \neq a^*} \inf_{\nu': \mu'_u \geq \mu'_{a^*}} \sum_{v \in N_{in}(u)} \omega_v G_{v,u} \text{KL}(\nu_u, \nu'_u) + \sum_{w \in N_{in}(a^*)} \omega_w G_{w,a^*} \text{KL}(\nu_{a^*}, \nu'_{a^*}), \\ &= \min_{u \neq a^*} \inf_{\nu': \mu'_u \geq \mu'_{a^*}} m_u \text{KL}(\nu_u, \nu'_u) + m_{a^*} \text{KL}(\nu_{a^*}, \nu'_{a^*}). \quad (m_u := \sum_{v \in N_{in}(u)} \omega_v G_{v,u}) \end{aligned}$$

Therefore, by **optimizing over  $\nu'$**  as in [GK16, Lemma 3] we obtain

$$\min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \geq \log(1/(2.4\delta)).$$

## Sample Complexity Lower Bounds - Uninformed Setting [Proof 3/3]

**Step 2 (Optimizing  $\nu'$ ):** We focus on **alternative models  $\nu'$**  that admit a different optimal vertex

$$\text{Alt}(\nu) = \cup_{v \neq a^*} \text{Alt}_v(\nu), \quad \text{Alt}_v(\nu) = \{\nu' \mid \mu'_v > \mu'_{a^*}\}.$$

Choose  $\nu'$  as to **minimize the LLR!**

$$\begin{aligned} & \inf_{\nu' \in \text{Alt}(\nu)} \sum_{u \in V} \sum_{v \in N_{in}(u)} \omega_v \text{KL}(\nu_{v,u}, \nu'_{v,u}) \\ &= \min_{u \neq a^*} \inf_{\nu': \mu'_u \geq \mu'_{a^*}} \sum_{v \in N_{in}(u)} \omega_v G_{v,u} \text{KL}(\nu_u, \nu'_u) + \sum_{w \in N_{in}(a^*)} \omega_w G_{w,a^*} \text{KL}(\nu_{a^*}, \nu'_{a^*}), \\ &= \min_{u \neq a^*} \inf_{\nu': \mu'_u \geq \mu'_{a^*}} m_u \text{KL}(\nu_u, \nu'_u) + m_{a^*} \text{KL}(\nu_{a^*}, \nu'_{a^*}). \quad (m_u := \sum_{v \in N_{in}(u)} \omega_v G_{v,u}) \end{aligned}$$

Therefore, by **optimizing over  $\nu'$**  as in [GK16, Lemma 3] we obtain

$$\min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \geq \log(1/(2.4\delta)).$$

# Sample Complexity Lower Bounds - Uninformed Setting - Bernoulli

- ▶ For **Bernoulli rewards** something funny happens in the **uninformed setting**<sup>3</sup>...
- ▶ Observing  $Z_{a,u} = 0$  can mean either “edge did not fire” or “reward was 0,”

$$P(Z = 0) = 1 - G_{v,u} \mu_u.$$

Because the learner never sees which edge fired, it is possible to construct an alternative model resemblingly perfectly the true model, under which an alternative arm is optimal!

## Proposition

*If  $(\nu_u)_{u \in V}$  are Bernoulli distributions with parameters  $(\mu_u)_{u \in V}$ , then  $a^*$  is unidentifiable, in the sense that  $(T^*(\nu))^{-1} = 0$ .*

---

<sup>3</sup>**Uninformed setting:** The learner does not know  $G$  nor which edge is activated at each time-step  $t$ .

# Sample Complexity Lower Bounds - Uninformed Setting - Bernoulli

- ▶ For **Bernoulli rewards** something funny happens in the **uninformed setting**<sup>3</sup>...
- ▶ Observing  $Z_{a,u} = 0$  can mean either “edge did not fire” or “reward was 0,”

$$P(Z = 0) = 1 - G_{v,u} \mu_u.$$

Because the learner never sees which edge fired, it is possible to construct an alternative model resemblingly perfectly the true model, under which an alternative arm is optimal!

## Proposition

*If  $(\nu_u)_{u \in V}$  are Bernoulli distributions with parameters  $(\mu_u)_{u \in V}$ , then  $a^*$  is unidentifiable, in the sense that  $(T^*(\nu))^{-1} = 0$ .*

---

<sup>3</sup>**Uninformed setting:** The learner does not know  $G$  nor which edge is activated at each time-step  $t$ .



# Sample Complexity Lower Bounds - Uninformed Setting - Bernoulli

- ▶ For **Bernoulli rewards** something funny happens in the **uninformed setting**<sup>3</sup>...
- ▶ Observing  $Z_{a,u} = 0$  can mean either “edge did not fire” or “reward was 0,”

$$P(Z = 0) = 1 - G_{v,u} \mu_u.$$

Because the learner never sees which edge fired, it is possible to construct an alternative model resemblingly perfectly the true model, under which an alternative arm is optimal!

## Proposition

*If  $(\nu_u)_{u \in V}$  are Bernoulli distributions with parameters  $(\mu_u)_{u \in V}$ , then  $a^*$  is unidentifiable, in the sense that  $(T^*(\nu))^{-1} = 0$ .*

---

<sup>3</sup>**Uninformed setting:** The learner does not know  $G$  nor which edge is activated at each time-step  $t$ .

# Sample Complexity Lower Bounds - Informed Setting - Bernoulli

What about the **informed setting**<sup>4</sup>?

- All good here, and the original sample complexity holds also for Bernoulli rewards.

$$\mathbb{E}_\nu[\tau] \geq T^*(\nu) \log \frac{1}{2.4\delta} \quad (2)$$

where

$$(T^*(\nu))^{-1} = \sup_{\omega \in \Delta(V)} \min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \text{ s.t. } m = G^\top \omega$$

---

<sup>4</sup>**Informed setting:** The learner either knows  $G$  or which edge was activated after choosing a node.

# Sample Complexity Lower Bounds - Informed Setting - Bernoulli

What about the **informed setting**<sup>4</sup>?

- All good here, and the original sample complexity holds also for Bernoulli rewards.

$$\mathbb{E}_\nu[\tau] \geq T^*(\nu) \log \frac{1}{2.4\delta} \quad (2)$$

where

$$(T^*(\nu))^{-1} = \sup_{\omega \in \Delta(V)} \min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \text{ s.t. } m = G^\top \omega$$

---

<sup>4</sup>**Informed setting:** The learner either knows  $G$  or which edge was activated after choosing a node.

## **TaS-FG: Track And Stop for Feedback Graphs**

---

# Components of a Strategy

A strategy is defined by

- ▶ Sampling rule
- ▶ Stopping rule
- ▶ Recommendation rule (we use the MLE)

# TaS-FG: Sampling Rule

How do we design an algorithm that **approaches the optimal sample complexity**?

$$T(\omega; \nu)^{-1} = \min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \text{ s.t. } m = G^\top \omega$$

The solution  $\omega^* \in \arg \inf_{\omega \in \Delta(V)} T(\omega; \nu)$  provides the best proportion of draws.

## Design

- ▶ Ensure that  $N_t/t$  (average selection frequency) tracks  $\omega^*(t)$  (computed w.r.t.  $\hat{\nu}(t)$ , the estimated model), where  $N_t$  is the visitation vector  $N(t) := [N_1(t) \ \dots \ N_K(t)]^\top$ .
- ▶ **Sampling rule:**

$$A_t \in \begin{cases} \arg \min_{u \in S_t} N_u(t) & \exists u : N_u(t) < \sqrt{t} - K/2 \\ \arg \min_{u \in V} N_u(t) - \sum_{n=1}^t \omega_u^*(n) & \text{otherwise} \end{cases}, \quad (3)$$

ensures  $\lim_{t \rightarrow \infty} \inf_{\omega \in C^*(\nu)} \|N(t)/t - \omega\|_\infty \rightarrow 0$  ( $C^*$  is the set of optimal allocations)<sup>5</sup>.

<sup>5</sup>Tracking a convex combination of all past solutions guarantees convergence to a unique point in  $C^*$ .

# TaS-FG: Sampling Rule

How do we design an algorithm that approaches the optimal sample complexity?

$$T(\omega; \nu)^{-1} = \min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \text{ s.t. } m = G^\top \omega$$

The solution  $\omega^* \in \arg \inf_{\omega \in \Delta(V)} T(\omega; \nu)$  provides the best proportion of draws.

## Design

► Ensure that  $N_t/t$  (average selection frequency) tracks  $\omega^*(t)$  (computed w.r.t.  $\hat{\nu}(t)$ , the estimated model), where  $N_t$  is the visitation vector  $N(t) := [N_1(t) \ \dots \ N_K(t)]^\top$ .

► Sampling rule:

$$A_t \in \begin{cases} \arg \min_{u \in S_t} N_u(t) & \exists u : N_u(t) < \sqrt{t} - K/2 \\ \arg \min_{u \in V} N_u(t) - \sum_{n=1}^t \omega_u^*(n) & \text{otherwise} \end{cases}, \quad (3)$$

ensures  $\lim_{t \rightarrow \infty} \inf_{\omega \in C^*(\nu)} \|N(t)/t - \omega\|_\infty \rightarrow 0$  ( $C^*$  is the set of optimal allocations)<sup>5</sup>.

<sup>5</sup>Tracking a convex combination of all past solutions guarantees convergence to a unique point in  $C^*$ .

# TaS-FG: Sampling Rule

How do we design an algorithm that **approaches the optimal sample complexity**?

$$T(\omega; \nu)^{-1} = \min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \text{ s.t. } m = G^\top \omega$$

The solution  $\omega^* \in \arg \inf_{\omega \in \Delta(V)} T(\omega; \nu)$  provides the best proportion of draws.

## Design

- Ensure that  $N_t/t$  (average selection frequency) tracks  $\omega^*(t)$  (computed w.r.t.  $\hat{\nu}(t)$ , the estimated model), where  $N_t$  is the visitation vector  $N(t) := [N_1(t) \ \dots \ N_K(t)]^\top$ .
- **Sampling rule:**

$$A_t \in \begin{cases} \arg \min_{u \in S_t} N_u(t) & \exists u : N_u(t) < \sqrt{t} - K/2 \\ \arg \min_{u \in V} N_u(t) - \sum_{n=1}^t \omega_u^*(n) & \text{otherwise} \end{cases}, \quad (3)$$

ensures  $\lim_{t \rightarrow \infty} \inf_{\omega \in C^*(\nu)} \|N(t)/t - \omega\|_\infty \rightarrow 0$  ( $C^*$  is the set of optimal allocations)<sup>5</sup>.

<sup>5</sup>Tracking a convex combination of all past solutions guarantees convergence to a unique point in  $C^*$ .



# TaS-FG: Stopping Rule

When do we stop?

$$T(\omega; \nu)^{-1} = \min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \text{ s.t. } m = G^\top \omega$$

## Stopping Rule

- ▶ The lower bound tells us that  $\tau \sim T^*(\nu) \log(1/\delta)$ . But we don't know the model! Additional price to pay  $O(\log \log(t))$ .
- ▶ Stopping as soon as

$$t \approx T(N(t)/t; \hat{\nu}(t))^6 \left[ \log \left( \frac{K-1}{\delta} \right) + O(\log \log(t)) \right]$$

guarantees correctness

$$\mathbb{P}_\nu(\tau < \infty, \hat{a}_\tau \neq a^*(\mu)) \leq \delta.$$

- ▶ With the previous sampling rule, we can guarantee  $\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_\nu[\tau]}{\ln(1/\delta)} \leq T^*(\nu)$ .

<sup>6</sup>Empirical characteristic time based on the MLE

# TaS-FG: Stopping Rule

When do we stop?

$$T(\omega; \nu)^{-1} = \min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \text{ s.t. } m = G^\top \omega$$

## Stopping Rule

- ▶ The lower bound tells us that  $\tau \sim T^*(\nu) \log(1/\delta)$ . But we don't know the model! Additional price to pay  $O(\log \log(t))$ .
- ▶ Stopping as soon as

$$t \approx T(N(t)/t; \hat{\nu}(t))^6 \left[ \log \left( \frac{K-1}{\delta} \right) + O(\log \log(t)) \right]$$

guarantees correctness

$$\mathbb{P}_\nu(\tau < \infty, \hat{a}_\tau \neq a^*(\mu)) \leq \delta.$$

- ▶ With the previous sampling rule, we can guarantee  $\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_\nu[\tau]}{\ln(1/\delta)} \leq T^*(\nu)$ .

<sup>6</sup>Empirical characteristic time based on the MLE

# TaS-FG: Stopping Rule

When do we stop?

$$T(\omega; \nu)^{-1} = \min_{u \neq a^*} (m_u + m_{a^*}) I_{\frac{m_{a^*}}{m_u + m_{a^*}}}(\nu_{a^*}, \nu_u) \text{ s.t. } m = G^\top \omega$$

## Stopping Rule

- ▶ The lower bound tells us that  $\tau \sim T^*(\nu) \log(1/\delta)$ . But we don't know the model! Additional price to pay  $O(\log \log(t))$ .
- ▶ Stopping as soon as

$$t \approx T(N(t)/t; \hat{\nu}(t))^6 \left[ \log \left( \frac{K-1}{\delta} \right) + O(\log \log(t)) \right]$$

guarantees correctness

$$\mathbb{P}_\nu(\tau < \infty, \hat{a}_\tau \neq a^*(\mu)) \leq \delta.$$

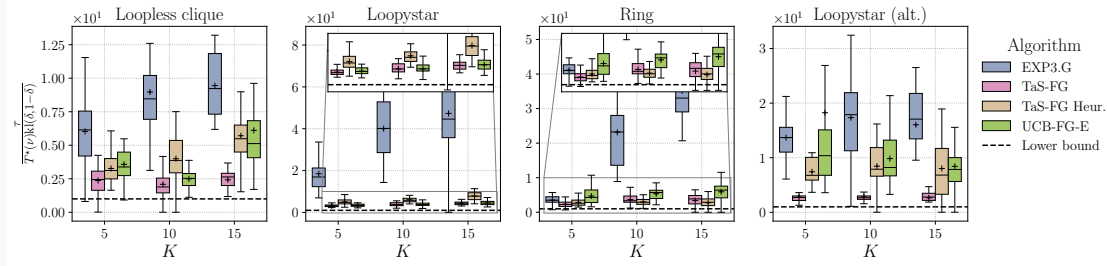
- ▶ With the previous sampling rule, we can guarantee  $\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_\nu[\tau]}{\ln(1/\delta)} \leq T^*(\nu)$ .

<sup>6</sup>Empirical characteristic time based on the MLE

## Numerical Results

---

# Numerical Results






Box plots of the normalized sample complexity  $\frac{\tau}{T^*(\nu)\text{kl}(\delta, 1-\delta)}$  for  $\delta = e^{-7}$  over 100 seeds. Boxes indicate the interquartile range, while the median and mean values are, respectively, the solid line and the  $+$  sign in black.

## Thank you for listening!

- Github repo:

<https://github.com/rssalessio/Pure-Exploration-with-Feedback-Graphs>



-  Noga Alon, Nicolo Cesa-Bianchi, Ofer Dekel, and Tomer Koren, *Online learning with feedback graphs: Beyond bandits*, Conference on Learning Theory, PMLR, 2015, pp. 23–35.
-  Aurélien Garivier and Emilie Kaufmann, *Optimal best arm identification with fixed confidence*, Conference on Learning Theory, PMLR, 2016, pp. 998–1027.
-  Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier, *On the complexity of best-arm identification in multi-armed bandit models*, The Journal of Machine Learning Research **17** (2016), no. 1, 1–42.