

# Structured Bandits

Baiyu Wang

University of Southampton

- 1 Introduction
  - Review of Multi-Armed and Contextual Bandits
  - Motivation
  - Structured Bandit
- 2 Optimism for Structured Bandits
  - UCB for Structured Bandit
  - Euler Dimension
  - Suboptimality of Optimism
- 3 Decision-Estimation Coefficient
  - E2D Algorithm
  - Decision-Estimation Coefficient
- 4 Conclusion

## Multi-Armed Bandit Protocol

**for**  $t = 1, \dots, T$  **do**

Select decision  $\pi^t \in \Pi := \{1, \dots, A\}$ .

Observe reward  $r^t \in \mathcal{R}$

## Contextual Bandit Protocol

**for**  $t = 1, \dots, T$  **do**

Observe context  $x^t \in \mathcal{X}$

Select decision  $\pi^t \in \Pi := \{1, \dots, A\}$ .

Observe reward  $r^t \in \mathcal{R}$

Measure performance for Bandit problem:

$$\mathbf{Reg} := \sum_{t=1}^T f^*(\pi^*) - \sum_{t=1}^T \mathbb{E}_{\pi^t \sim p^t} [f^*(\pi^t)]$$

# Practical Issue

- In medicine, the treatment may be a continuous variable, such as a dosage. The treatment could even be a high-dimensional vector (such as dosages for many different medications).
- In pricing applications, a seller might aim to select a continuous price or vector in order to maximize their returns.
- In routing applications, the decision space may be finite, but combinatorially large. For example, the decision might be a path or flow in a graph.

⇒ Decision space  $\Pi$  is not finite set.

# Regret Bound

## Multi-Armed Bandit:

- $\varepsilon$ -Greedy algorithm:  $\mathbf{Reg} \lesssim A^{1/3} T^{2/3} \cdot \log^{1/3}(AT/\delta)$ .
- UCB algorithm:  $\mathbf{Reg} \lesssim \sqrt{AT \log(AT/\delta)}$ .
- Posterior Sampling Algorithm:  $\mathbf{Reg} \lesssim \sqrt{AT \log(A)} / \sqrt{AT \log |\mathcal{F}|}$
- Exp3 Algorithm:  $\mathbf{Reg} \lesssim \sqrt{AT \log A}$

## Contextual Bandit:

- $\varepsilon$ -Greedy algorithm:  $\mathbf{Reg} \lesssim A^{1/3} T^{2/3} \cdot \mathbf{Est}_{\text{Sq}}(\mathcal{F}, T, \delta)^{1/3}$ .
- LinUCB algorithm:  $\mathbf{Reg} \lesssim \sqrt{dT \log(|\mathcal{F}|/\delta) \log(1 + T/d)}$

# Necessity of structural assumptions

Let  $\Pi = [A]$ , and let  $\mathcal{F} = \{f_i\}_{i \in [A]}$ , where

$$f_i(\pi) := \frac{1}{2} + \frac{1}{2} \mathbb{I}\{\pi = i\}.$$

- $\text{Reg} \gtrsim A$  and  $\log |\mathcal{F}| = \log(A)$  for this setting.
- $\text{Reg} \lesssim \sqrt{AT \log |\mathcal{F}|}$  is impossible if  $A \gg T$ .

# Structured Bandit

## Structured Bandit Protocol

**for**  $t = 1, \dots, T$  **do**

Select decision  $\pi^t \in \Pi$ .  *$\Pi$  is large and potentially continuous.*

Observe reward  $r^t \in \mathbb{R}$ .



# Assumption and Regret

- Decision:  $\pi^t \in \Pi$ , for  $t = 1, \dots, T$ .
- Stochastic rewards: Rewards are generated independently via  $r^t \sim M^*(\cdot|\pi^t)$ .
- Mean reward function:  $f^*(\pi) := \mathbb{E}[r|\pi]$ .
- Regret function:  $\mathbf{Reg} := \sum_{t=1}^T f^*(\pi^*) - \sum_{t=1}^T \mathbb{E}_{\pi^t \sim p^t} [f^*(\pi^t)]$ ,  
where  $\pi^* := \arg \max_{\pi \in \Pi} f^*(\pi)$ .
- Data history:  $\mathcal{H}^t = (\pi^1, r^1), \dots, (r^t, \pi^t)$ .
- The decision-maker has access to a class  $\mathcal{F} \subset \{f: \Pi \rightarrow \mathbb{R}\}$  such that  $f^* \in \mathcal{F}$ .



# Upper Confidence Bound

## Assumption and definition:

- Assume  $\mathcal{F} = \{f: \Pi \rightarrow [0, 1]\}$  and  $r^t \in [0, 1]$  almost surely.
- $\hat{f}^t = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^{t-1} (f(\pi^i) - r^i)^2$ .
- Confidence sets  $\mathcal{F}^1 = \mathcal{F}$  and

$$\mathcal{F}^t = \left\{ f \in \mathcal{F} : \sum_{i=1}^{t-1} (f(\pi^i) - r^i)^2 \leq \sum_{i=1}^{t-1} (\hat{f}^t(\pi^i) - r^i)^2 + \beta \right\},$$

where  $\beta := 8 \log(|\mathcal{F}|/\delta)$ .

- $\bar{f}^t(\pi) := \max_{f \in \mathcal{F}^t} f(\pi)$  is the upper confidence bound.

## Generalized UCB algorithm:

$$\pi^t = \arg \max_{\pi \in \Pi} \bar{f}^t(\pi)$$

# Upper Confidence Width

## Lemma 10

Let  $\pi^1, \dots, \pi^T$  be chosen by an arbitrary (and possibly randomized) decision-making algorithm. With probability at least  $1 - \delta$ ,  $f^* \in \mathcal{F}^t$  for all  $t \in [T]$ . Moreover, with probability at least  $1 - \delta$ , for all  $\tau \leq T$ , all  $f \in \mathcal{F}^\tau$  with  $\beta = 8 \log(|\mathcal{F}/\delta|)$  satisfy

$$\sum_{t=1}^{\tau} \mathbb{E}_{\pi^t \sim p^t} [(f(x^t, \pi^t) - f^*(x^t, \pi^t))^2] \leq 4\beta.$$

From [Lemma 10](#),  $f^* \in \mathcal{F}$  with high probability. The regret is bounded by upper confidence width

$$\text{Reg} \leq \sum_{t=1}^T \left( \bar{f}^t(\pi^t) - f^*(\pi^t) \right).$$

When will the confidence shrink?

# LinearUCB Algorithm

We hope  $\Pi$  and  $\mathcal{F}$  have nice structure to get a better regret bound.

Recall linearUCB algorithm:

$$\mathcal{F} = \{ \pi \mapsto \langle \theta, \phi(\pi) \rangle \mid \theta \in \Theta \subset \mathbb{B}_2^d(1) \}$$

From [Proposition 7](#) we know  $\mathbf{Reg} \lesssim \sqrt{dT \log |\mathcal{F}|}$ .

Is there a general version when we move beyond linear model?

# Definition

## Eulder Dimension

$\mathcal{F} \subset (\Pi \rightarrow \mathbb{R})$  and  $f^* : \Pi \rightarrow \mathbb{R}$  be given, and define  $\underline{\text{Edim}}_{f^*}(\mathcal{F}, \varepsilon)$  as the length of the longest sequence of decisions  $\pi^1, \dots, \pi^d \in \Pi$  such that for all  $t \in [d]$ , there exists  $f^t \in \mathcal{F}$  such that

$$|f^t(\pi^t) - f^*(\pi^t)| > \varepsilon, \quad \text{and} \quad \sum_{i < t} (f^i(\pi^i) - f^*(\pi^i))^2 \leq \varepsilon^2$$

Eluder dimension:  $\text{Edim}_{f^*}(\mathcal{F}, \varepsilon) = \sup_{\varepsilon' \geq \varepsilon} \underline{\text{Edim}}_{f^*}(\mathcal{F}, \varepsilon') \vee 1$ ,  
 $\text{Edim}(\mathcal{F}, \varepsilon) = \max_{f^* \in \mathcal{F}} \text{Edim}_{f^*}(\mathcal{F}, \varepsilon).$

- When  $\beta = \varepsilon^2$  in confidence set, then

$$\sum_{t=1}^T \mathbb{I} \{ \bar{f}^t(\pi^t) - f^*(\pi^t) > \varepsilon \} \leq \text{Edim}_{f^*}(\mathcal{F}, \varepsilon).$$

- Monotonicity with respect to  $\varepsilon$ .

# Regret Bound for UCB

## Regret Bound

For a finite set of functions  $\mathcal{F} \subset (\Pi \rightarrow [0, 1])$ , using  $\beta = 8 \log(|\mathcal{F}|/\delta)$ , the generalized UCB algorithm guarantees that with probability at least  $1 - \delta$ ,

$$\begin{aligned} \mathbf{Reg} &\lesssim \min_{\varepsilon > 0} \{ \sqrt{\text{Edim}(\mathcal{F}, \varepsilon) \cdot T \log(|\mathcal{F}|/\delta)} + \varepsilon T \} \\ &\lesssim \sqrt{\text{Edim}(\mathcal{F}, T^{-1/2}) \cdot T \log(|\mathcal{F}|/\delta)} \end{aligned}$$

# Proof of Regret Bound

Define

$$\overline{\mathcal{F}}^t = \left\{ f \in \mathcal{F} \mid \sum_{i < t} (f(\pi^i) - f^*(\pi^i))^2 \leq 4\beta \right\}.$$

By [Lemma 10](#),  $f^* \in \mathcal{F}^t$  and  $\mathcal{F}^t \subseteq \overline{\mathcal{F}}^t$  for  $t = 1, \dots, T$  with probability at least  $1 - \delta$ .

Define

$$w^t(\pi) = \sup_{f \in \overline{\mathcal{F}}^t} [f(\pi) - f^*(\pi)].$$

The regret can be bounded by

$$\mathbf{Reg} \leq \sum_{t=1}^T \bar{f}^t(\pi^t) - f^*(\pi^t) \leq \sum_{t=1}^T w^t(\pi^t)$$

Fix  $\varepsilon > 0$ ,

$$\sum_{t=1}^T w^t(\pi^t) \leq \sum_{t=1}^T w^t(\pi^t) \mathbb{I}\{w^t(\pi^t) > \varepsilon\} + \varepsilon T.$$

### Lemma 1

Fix a function class  $\mathcal{F}$ , function  $f^* \in \mathcal{F}$ , and parameter  $\beta > 0$ . For any sequence  $\pi^1, \dots, \pi^T$ , if we define

$$w^t(\pi) = \sup_{f \in \mathcal{F}} \left\{ f(\pi) - f^*(\pi) : \sum_{i < t} (f(\pi^i) - f^*(\pi^i))^2 \leq \beta \right\}$$

then for all  $\alpha > 0$ ,

$$\sum_{t=1}^T \mathbb{I}\{w^t(\pi^t) > \alpha\} \leq \left( \frac{\beta}{\alpha^2} + 1 \right) \cdot \underline{\text{Edim}}_{f^*}(\mathcal{F}, \alpha)$$

Order the indices  $\{1, \dots, T\}$  as  $\{i_1, \dots, i_T\}$ , so that

$$w^{i_1}(\pi^{i_1}) \geq w^{i_2}(\pi^{i_2}) \geq \dots \geq w^{i_T}(\pi^{i_T}).$$

Consider any index  $\tau$  for which  $w^{i_\tau}(\pi^{i_\tau}) > \varepsilon$ . For any  $\alpha > \varepsilon$ , if we have  $w^{i_\tau}(\pi^{i_\tau}) > \alpha$ , [Lemma 1](#) implies that

$$\tau \leq \sum_{t=1}^T \mathbb{I}\{w^t(\pi^t) > \alpha\} \leq \left(\frac{4\beta}{\alpha^2} + 1\right) \underline{\text{Edim}}_{\mathcal{F}}(\mathcal{F}, \alpha) \leq \frac{5\beta}{\alpha^2} \underline{\text{Edim}}_{\mathcal{F}}(\mathcal{F}, \alpha).$$

Hence,

$$w^{i_\tau}(\pi^{i_\tau}) \leq \sqrt{\frac{5\beta \underline{\text{Edim}}(\mathcal{F}, \varepsilon)}{\tau}}$$





# Proof of Lemma 1

The proof will be finished by following steps:

1. Define  $\alpha$ -dependent, which is similar to the Eulder dimension.
2. For  $w^t(\pi^t) > \alpha$  in a sequence  $\pi^1, \dots, \pi^t$ , we find at most  $\frac{\beta}{\alpha^2}$  of disjoint subsequences that are  $\alpha$ -dependent of  $\pi^t$ .
3. For any sequence, we find that there exists one  $\pi$  in the sequence  $\alpha$ -dependent with at least  $\lfloor \tau/d \rfloor$  disjoint subsequence.
4. Establish relationship between Step 2 and Step 3.

## Step 1:

Let  $d = \underline{\text{Edim}}_{f^*}(\mathcal{F}, \alpha)$ .

$\alpha$ -independent:

$\pi$  is  $\alpha$ -independent of  $\pi^1, \dots, \pi^t$  if  $\exists f \in \mathcal{F}$   
 such that  $|f(\pi) - f^*(\pi)| > \alpha$  and  $\sum_{i=1}^t (f(\pi^i) - f^*(\pi^i))^2 \leq \alpha^2$ .

$\alpha$ -dependent:

$\pi$  is  $\alpha$ -dependent on  $\pi^1, \dots, \pi^t$  if  
 $\forall f \in \mathcal{F}$  with  $\sum_{i=1}^t (f(\pi^i) - f^*(\pi^i))^2 \leq \alpha^2, |f(\pi) - f^*(\pi)| \leq \alpha$ .

## Step 2:

For any  $t$ , if  $w^t(\pi^t) > \alpha$ , then  $\pi_t$  is  $\alpha$ -dependent on at most  $\beta/\alpha^2$  disjoint subsequences of  $\pi^1, \dots, \pi^{t-1}$ .

Let  $f$  be  $|f(\pi^t) - f^*(\pi^t)| > \alpha$ . If  $\pi^t$  is  $\alpha$ -dependent on  $\pi^{i_1}, \dots, \pi^{i_k}$  but  $w^t(\pi^t) > \alpha$ , we must have

$$\sum_{j=1}^k (f(\pi^{i_j}) - f^*(\pi^{i_j}))^2 \geq \alpha^2$$

If there are  $M$  such disjoint sequences, we have

$$M\alpha^2 \leq \sum_{i < t} (f(\pi^i) - f^*(\pi^i))^2 \leq \beta$$

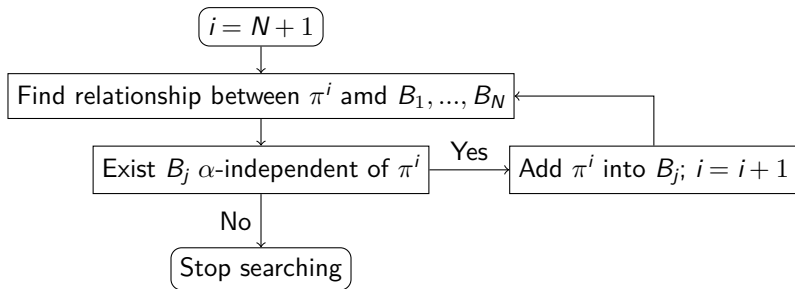
so  $M \leq \frac{\beta}{\alpha^2}$ .

### Step 3:

For  $\tau$  and any sequence  $(\pi^1, \dots, \pi^\tau)$ , there is some  $\pi^j$  is  $\alpha$ -dependent on at least  $\lfloor \tau/d \rfloor$  disjoint subsequences of  $\pi^1, \dots, \pi^{j-1}$ .

Let  $N = \lfloor \tau/d \rfloor$ , and let  $B_1, \dots, B_N$  be subsequences of  $\pi^1, \dots, \pi^\tau$ , and  $B_i = (\pi^i)$ .

Do the following step



If stop at  $i = \tau$ , we have  $\sum_{i=1}^N |B_i| \geq \tau \geq dN$ . Moreover,  $|B_i| \leq d$ , so in this case  $\pi^\tau$  is  $\alpha$ -dependent on all  $B_i$ .

## Step 4:

Let  $(\pi^{t_1}, \dots, \pi^{t_\tau})$  be the subsequence  $\pi^1, \dots, \pi^T$  consisting of all elements for which  $w^{t_i}(\pi^{t_i}) > \alpha$ .

From **Step 2**, each  $\pi^{t_i}$  is dependent on at most  $\beta/\alpha^2$  disjoint subsequences of  $(\pi^{t_1}, \dots, \pi^{t_\tau})$ .

From **Step 3**, one element is dependent on at least  $\lfloor \tau/d \rfloor$  disjoint subsequences.

We must have  $\lfloor \tau/d \rfloor \leq \beta/\alpha^2$ , and which implies that  $\tau \leq (\beta/\alpha^2 + 1) d$ .

## When will this bound be vacuous?

Consider generalised linear models

$$\mathcal{F} = \{\pi \mapsto \sigma(\langle \theta, \phi(\pi) \rangle) \mid \theta \in \Theta \subset \mathbb{B}_2^d(1)\}$$

Let  $\sigma(z) = +\text{relu}(z)$  or  $\sigma(z) = -\text{relu}(z)$ , where  $\text{relu}(z) := \max\{z, 0\}$  is the ReLU function.

The lower bound of Euler dimension is

$$\text{Edim}(\mathcal{F}, \varepsilon) \gtrsim e^d$$

for constant  $\varepsilon$ .

# Suboptimality of Optimism

The UCB algorithm is useful for some special cases, it does not attain optimal regret for any structured bandit problem. We see an example by adding "cheating" action into decision space.

**Example:** Let  $A \in \mathbb{N}$  be a power of 2 and consider the following  $\mathcal{F}$ .

- The decision space is  $\Pi = [A] \cup \mathcal{C}$ , where  $\mathcal{C} = \{c_1, \dots, c_{\log_2(A)}\}$  is a set of "cheating" actions.
- For all actions  $\pi \in [A]$ ,  $f(\pi) \in [0, 1]$  for all  $f \in \mathcal{F}$ .
- For each  $f \in \mathcal{F}$ , rewards for actions in  $\mathcal{C}$  take the following form. Let  $b(f) = (b_1(f), \dots, b_{\log_2(A)}(f)) \in \{0, 1\}^{\log_2(A)}$  be a binary encoding for the index of  $\pi_f \in [A]$ . For each action  $c_i \in \mathcal{C}$ , we set

$$f(c_i) = -b_i(f).$$



## Simple analyses

- Suppose that rewards are Gaussian with  $r \sim \mathcal{N}(f^*(\pi), 1)$  under  $\pi$ .
- If we do this for each  $c_i \in \mathcal{C}$ , which will incur  $\tilde{O}(\log_2(A))$  regret.
- We can stop exploring, and commit to playing  $\pi_{f^*}$  for remaining rounds.

$\Rightarrow$  with probability at least  $1 - \delta$ ,

$$\mathbf{Reg} \lesssim \log_2^2(A/\delta).$$

## UCB algorithm

- $c_i \in \mathcal{C}$  have  $f(c_i) \leq 0$  for all  $f \in \mathcal{F}$ , we have  $\bar{f}^t(c_i) \leq 0$ .
- UCB algorithm only pull actions in  $[A]$ , ignoring the cheating actions.

$\Rightarrow$

$$\mathbf{Reg} \gtrsim \sqrt{AT}.$$

**The "cheating" actions are guaranteed to have low reward, UCB avoids them even though they are very informative.**

We conclude that:

- Balance the tradeoff between optimizing reward and acquiring information to get optimal regret bound.
- The Inverse Gap Weighting algorithm attained optimal sample complexity for any choice of class  $\mathcal{F}$ , and all that needed to change was how to perform estimation.

# Online Regression Oracle

## Online Regression Oracle

At each time  $t \in [T]$ , an online regression oracle returns, given

$$(\pi^1, r^1), \dots, (\pi^{t-1}, r^{t-1})$$

with  $\mathbb{E}[r^i | \pi^i] = f^*(\pi^i)$  and  $\pi^i \sim p^i$ , a function  $\hat{f}^t : \Pi \rightarrow \mathbb{R}$  such that

$$\sum_{t=1}^T \mathbb{E}_{\pi^t \sim p^t} \left( \hat{f}^t(\pi^t) - f^*(\pi^t) \right)^2 \leq \text{Est}_{\text{Sq}}(\mathcal{F}, T, \delta)$$

with probability at least  $1 - \delta$ . Here,  $p^i(\cdot | \mathcal{H}^{i-1})$  is the randomization distribution for the decision-maker.

# E2D Algorithm

## Estimation-to-Decision (E2D) for Structured Bandits

Input: Exploration parameter  $\gamma > 0$ .

for  $t = 1, \dots, T$  do

Obtain  $\hat{f}^t$  from online regression oracle with  $(\pi^1, r^1), \dots, (\pi^{t-1}, r^{t-1})$ .

Compute

$$p^t = \arg \min_{p \in \Delta(\Pi)} \max_{f \in \mathcal{F}} \mathbb{E}_{\pi \sim p} \left[ f(\pi_f) - f(\pi) - \gamma \cdot \left( f(\pi) - \hat{f}^t(\pi) \right)^2 \right].$$

Select action  $\pi^t \sim p^t$ .

# Decision-to-Estimation Coefficient

**Decision-Estimation Coefficient** is a complexity measure for  $\mathcal{F}$ , and it is defined as:

$$\text{dec}_\gamma(\mathcal{F}, \hat{f}) = \min_{p \in \Delta(\Pi)} \max_{f \in \mathcal{F}} \mathbb{E}_{\pi \sim p} \left[ \underbrace{f(\pi_f) - f(\pi)}_{\text{regret of decision}} - \gamma \cdot \underbrace{(f(\pi) - \hat{f}(\pi))^2}_{\text{information gain for obs.}} \right]$$

$$\text{dec}_\gamma(\mathcal{F}) = \sup_{\hat{f} \in \text{co}(\mathcal{F})} \text{dec}_\gamma(\mathcal{F}, \hat{f})$$

# Regret Bound for E2D

## Regret Bound

The E2D algorithm with exploration parameter  $\gamma > 0$  guarantees that with probability at least  $1 - \delta$ ,

$$\mathbf{Reg} \leq \text{dec}_\gamma(\mathcal{F}) \cdot T + \gamma \cdot \text{Est}_{\text{sq}}(\mathcal{F}, T, \delta)$$

We can optimize over the parameter  $\gamma$  in the result above, which yields

$$\begin{aligned} \mathbf{Reg} &\leq \inf_{\gamma > 0} \{ \text{dec}_\gamma(\mathcal{F}) \cdot T + \gamma \cdot \text{Est}_{\text{sq}}(\mathcal{F}, T, \delta) \} \\ &\leq 2 \cdot \inf_{\gamma > 0} \max \{ \text{dec}_\gamma(\mathcal{F}) \cdot T, \gamma \cdot \text{Est}_{\text{sq}}(\mathcal{F}, T, \delta) \} \end{aligned}$$

For finite classes, the exponential weights method give  $\text{Est}_{\text{sq}}(\mathcal{F}, T, \delta) \lesssim \log(|\mathcal{F}|/\delta)$ , and this bound specializes to

$$\mathbf{Reg} \lesssim \inf_{\gamma > 0} \max \{ \text{dec}_\gamma(\mathcal{F}) \cdot T, \gamma \cdot \log(|\mathcal{F}|/\delta) \}$$

# Proof of Regret Bound

We write

$$\begin{aligned}
 \mathbf{Reg} &= \sum_{t=1}^T \mathbb{E}_{\pi^t \sim p^t} [f^*(\pi^*) - f^*(\pi^t)] \\
 &= \sum_{t=1}^T \mathbb{E}_{\pi^t \sim p^t} [f^*(\pi^*) - f^*(\pi^t)] - \gamma \cdot \mathbb{E}_{\pi^t \sim p^t} \left[ \left( f^*(\pi^t) - \hat{f}^t(\pi^t) \right)^2 \right] \\
 &\quad + \gamma \cdot \mathbf{Est}_{\text{Sq}}(\mathcal{F}, T, \delta)
 \end{aligned}$$

For each  $t$ , since  $f^* \in \mathcal{F}$ , we have

$$\begin{aligned}
 & \mathbb{E}_{\pi^t \sim p^t} [f^*(\pi^*) - f^*(\pi^t)] - \gamma \cdot \mathbb{E}_{\pi^t \sim p^t} \left[ \left( f^*(\pi^t) - \hat{f}^t(\pi^t) \right)^2 \right] \\
 & \leq \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}_{\pi^t \sim p^t} [f(\pi_f) - f(\pi^t)] - \gamma \cdot \mathbb{E}_{\pi^t \sim p^t} \left[ \left( f(\pi^t) - \hat{f}^t(\pi^t) \right)^2 \right] \right\} \\
 & = \inf_{p \in \Delta(\Pi)} \sup_{f \in \mathcal{F}} \mathbb{E}_{\pi \sim p} \left[ f(\pi_f) - f(\pi) - \gamma \cdot \left( f(\pi^t) - \hat{f}^t(\pi^t) \right)^2 \right] \\
 & = \text{dec}_\gamma \left( \mathcal{F}, \hat{f}^t \right)
 \end{aligned}$$

We conclude that

$$\mathbf{Reg} \leq \sup_{\hat{f}} \text{dec}_\gamma(\mathcal{F}, \hat{f}) \cdot T + \gamma \cdot \mathbf{Est}_{\text{Sq}}(\mathcal{F}, T, \delta)$$



# Example: Multi-Armed Bandit

## Decision-to-Estimation Bound (IGW minimizes the DEC)

For the Multi-Armed Bandit setting, where  $\Pi = [A]$  and  $\mathcal{F} = \mathbb{R}^A$

- the Inverse Gap Weighting distribution  $p = \text{IGW}_{4\gamma}(\hat{f})$  is the exact minimizer for  $\text{dec}_\gamma(\mathcal{F}, \hat{f})$ .
- $\text{dec}_\gamma(\mathcal{F}, \hat{f}) = \frac{A+1}{4\gamma}$ .

## Inverse Gap Weighting

Give a vector  $\hat{f} = (\hat{f}(1), \dots, \hat{f}(A)) \in \mathbb{R}^A$ , the Inverse Gap Weighting distribution  $p = \text{IGW}_\gamma(\hat{f}(1), \dots, \hat{f}(A))$  with parameter  $\gamma \geq 0$  is defined as

$$p(\pi) = \frac{1}{\lambda + 2\gamma(\hat{f}(\hat{\pi}) - \hat{f}(\pi))},$$

where  $\hat{\pi} = \arg \max_{\pi} \hat{f}(\pi)$ , and  $\lambda \in [1, A]$  satisfy  $\sum_{\pi} p(\pi) = 1$

## Proposition 9

Consider a finite decision space  $\Pi = \{1, \dots, A\}$ . For any vector  $\hat{f} \in \mathbb{R}^A$  and  $\gamma > 0$ , define  $p = \text{IGW}_\gamma(\hat{f}(1), \dots, \hat{f}(A))$ . This strategy guarantees that for all  $f^* \in \mathbb{R}^A$ ,

$$\mathbb{E}[f^*(\pi^*) - f^*(\pi)] \leq \frac{A}{\gamma} + \gamma \cdot \mathbb{E}_{\pi \sim p}[(\hat{f}(\pi) - f^*(\pi))^2].$$

- By rewriting [Proposition 9](#), we can deduce that the DEC is bounded by  $\frac{A+1}{4\gamma}$ .
- [Proposition 14](#) shows that IGW is actually the best possible distribution for this minimax problem.

# Proof of DEC Bound

Rewrite the definition of DEC:

$$\begin{aligned}
 & \min_{p \in \Delta([A])} \max_{f \in \mathbb{R}^A} \mathbb{E}_{\pi \sim p} \left[ f(\pi_f) - f(\pi) - \gamma(f(\pi) - \hat{f}(\pi))^2 \right] \\
 &= \min_{p \in \Delta([A])} \max_{f \in \mathbb{R}^A} \max_{\pi^* \in [A]} \mathbb{E}_{\pi \sim p} \left[ f(\pi^*) - f(\pi) - \gamma(f(\pi) - \hat{f}(\pi))^2 \right] \\
 &= \min_{p \in \Delta([A])} \max_{\pi^* \in [A]} \max_{f \in \mathbb{R}^A} \mathbb{E}_{\pi \sim p} \left[ f(\pi^*) - f(\pi) - \gamma(f(\pi) - \hat{f}(\pi))^2 \right].
 \end{aligned}$$

Proof step:

- IGW achieve lower bound of

$$\text{Inner Equaiton} = \max_{\pi^* \in [A]} \max_{f \in \mathbb{R}^A} \mathbb{E}_{\pi \sim p} \left[ f(\pi^*) - f(\pi) - \gamma(f(\pi) - \hat{f}(\pi))^2 \right]$$

- IGW achieve the condition of

$$\min_{p \in \Delta([A])} \text{Inner Equaiton}$$

# Step 1

**Firstly, we find the lower bound.**

For any fixed  $p$  and  $\pi^*$ , first-order conditions for optimality imply that

$$f(\pi) = \hat{f}(\pi) - \frac{1}{2\gamma} + \frac{1}{2\gamma p(\pi^*)} \mathbb{I}\{\pi = \pi^*\}$$

It gives that

$$\mathbb{E}_{\pi \sim p} [f(\pi^*) - f(\pi)] = \mathbb{E}_{\pi \sim p} [\hat{f}(\pi^*) - \hat{f}(\pi)] + \frac{1 - p(\pi^*)}{2\gamma p(\pi^*)}$$

and

$$\gamma \mathbb{E}_{\pi \sim p} [(f(\pi) - \hat{f}(\pi))^2] = \frac{1 - p(\pi^*)}{4\gamma} + \frac{(1 - p(\pi^*))^2}{4\gamma p(\pi^*)} = \frac{1}{4\gamma p(\pi^*)} - \frac{1}{4\gamma}$$

For any  $p \in \Delta(\Pi)$ ,

$$\begin{aligned}
 & \max_{\pi^* \in [A]} \max_{f \in \mathbb{R}^A} \mathbb{E}_{\pi \sim p} \left[ f(\pi^*) - f(\pi) - \gamma(f(\pi) - \hat{f}(\pi))^2 \right] \\
 &= \max_{\pi^* \in [A]} \left\{ \mathbb{E}_{\pi \sim p} \left[ \hat{f}(\pi^*) - \hat{f}(\pi) \right] + \frac{1}{4\gamma p(\pi^*)} \right\} - \frac{1}{4\gamma} \\
 &\geq \mathbb{E}_{\pi^* \sim p} \left[ \mathbb{E}_{\pi \sim p} \left[ \hat{f}(\pi^*) - \hat{f}(\pi) \right] + \frac{1}{4\gamma p(\pi^*)} \right] - \frac{1}{4\gamma} \\
 &= \frac{A}{4\gamma} - \frac{1}{4\gamma}
 \end{aligned}$$

**Next, we prove that IGW can achieve this lower bound.**

Let  $p = \text{IGW}_{4\gamma}(\hat{f})$ , for all  $\pi^*$ ,

$$\begin{aligned} & \mathbb{E}_{\pi \sim p} [\hat{f}(\pi^*) - \hat{f}(\pi)] + \frac{1}{4\gamma p(\pi^*)} \\ &= \mathbb{E}_{\pi \sim p} [\hat{f}(\pi^*) - \hat{f}(\pi)] + \frac{\lambda}{4\gamma} + \hat{f}(\hat{\pi}) - \hat{f}(\pi^*) = \mathbb{E}_{\pi \sim p} [\hat{f}(\hat{\pi}) - \hat{f}(\pi)] + \frac{\lambda}{4\gamma} \end{aligned}$$

Hence,

$$\begin{aligned} & \max_{\pi^* \in [A]} \left\{ \mathbb{E}_{\pi \sim p} [\hat{f}(\pi^*) - \hat{f}(\pi)] + \frac{1}{4\gamma p(\pi^*)} \right\} \\ &= \min_{\pi^* \in [A]} \left\{ \mathbb{E}_{\pi \sim p} [\hat{f}(\pi^*) - \hat{f}(\pi)] + \frac{1}{4\gamma p(\pi^*)} \right\} \\ &= \mathbb{E}_{\pi^* \sim p} \left\{ \mathbb{E}_{\pi \sim p} [\hat{f}(\pi^*) - \hat{f}(\pi)] + \frac{1}{4\gamma p(\pi^*)} \right\} = \frac{A}{4\gamma} \end{aligned}$$

$p = \text{IGW}_{4\gamma}(\hat{f})$  achieves the optimal value.

# Step 2

Relaxing to  $p \in \mathbb{R}_+^A$ . Define

$$g_{\pi^*}(p) = \hat{f}(\pi^*) + \frac{1}{4\gamma p(\pi^*)}.$$

Let  $\nu \in \mathbb{R}$  be a Lagrange multiplier, and Lagrangian is

$$\mathcal{L}(p, \nu) = g_{\pi^*}(p) - \sum_{\pi} p(\pi) \hat{f}(\pi) + \nu \left( \sum_{\pi} p(\pi) - 1 \right)$$

Based on **Step 1**, we aim to show that  $p \in \Delta(\Pi)$  is optimal for

$$\min_{p \in \Delta(\Pi)} \max_{\pi^* \in [A]} \left\{ \mathbb{E}_{\pi \sim p} [\hat{f}(\pi^*) - \hat{f}(\pi)] + \frac{1}{4\gamma p(\pi^*)} \right\} - \frac{1}{4\gamma}$$

It is sufficient to find  $\nu$  such that

$$0 \in \partial_p \mathcal{L}(p, \nu)$$

where  $\partial_p$  denotes the subgradient with respect to  $p$ .

For a convex function  $h(x) = \max_y g(x, y)$ , we have

$$\partial_x h(x) = \text{co} \left( \left\{ \nabla g(x, y) \mid g(x, y) = \max_{y'} g(x, y') \right\} \right).$$

Hence,

$$\partial_p \mathcal{L}(p, \nu) = \nu \mathbf{1} - \hat{f} + \text{co} \left( \left\{ \nabla_p g_{\pi^*}(p) \mid g_{\pi^*}(p) = \max_{\pi'} g_{\pi'}(p) \right\} \right)$$



Now, let  $p = \text{IGW}_{4\gamma}(\hat{f})$ . From **Step 1**, we know that  $g_\pi(p) = g_{\pi'}(p)$  for all  $\pi, \pi'$ , so

$$\partial_p \mathcal{L}(p, \nu) = \nu \mathbf{1} - \hat{f} + \text{co}(\{\nabla_p g_{\pi^*}(p)\}_{\pi^* \in \Pi}).$$

Since  $\nabla_p g_{\pi^*}(p) = -\frac{1}{4\gamma p^2(\pi^*)} \mathbf{e}_{\pi^*}$ ,

$$\delta := \sum_{\pi} p(\pi) \nabla_p g_{\pi}(p) = \left\{ -\frac{1}{4\gamma p(\pi)} \right\}_{\pi \in \Pi} = \left\{ -\frac{\lambda}{4\gamma} - \hat{f}(\hat{\pi}) + \hat{f}(\pi) \right\}_{\pi \in \Pi}.$$

Therefore,  $\delta \in \text{co}(\{\nabla_p g_{\pi^*}(p)\}_{\pi^* \in \Pi})$ , by choosing  $\nu = \frac{\lambda}{4\gamma} + \hat{f}(\hat{\pi})$ , we have

$$\nu \mathbf{1} - \hat{f} + \delta = 0,$$

which finishes the proof.

# Conclusion

- Introduce Structured Bandit, and generalise the decision space  $\Pi$  into large and potentially continuous space.
- Applying UCB algorithm with Euler dimension to solve Structured Bandit problem, where

$$\mathbf{Reg} \lesssim \min_{\varepsilon > 0} \{ \sqrt{\text{Edim}(\mathcal{F}, \varepsilon) \cdot T \log(|\mathcal{F}|/\delta)} + \varepsilon T \}.$$

- Using Estimation-to-Decision (E2D) for Structured Bandit, which provide the optimal rate. The regret bound is

$$\mathbf{Reg} \leq \text{dec}_{\gamma}(\mathcal{F}) \cdot T + \gamma \cdot \text{Est}_{\text{sq}}(\mathcal{F}, T, \delta)$$