

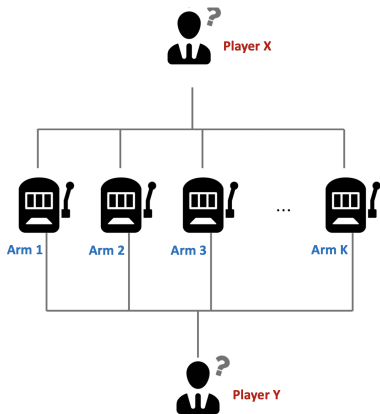
The Pareto Frontier of Instance-Dependent Guarantees in Multi-Player Multi-Armed Bandits with no Communication

Mark Sellke (Stanford)

With Allen Liu (MIT)

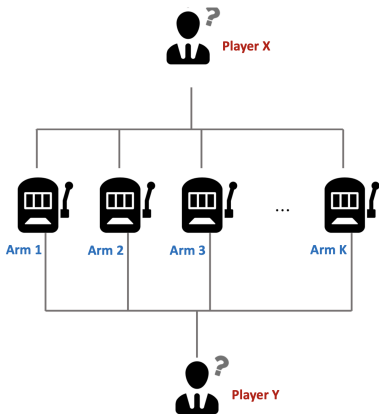


Multi-Player (Cooperative) Bandits



Consider $m > 1$ players with K stochastic bandit arms. Each time $1 \leq t \leq T$, each player X chooses an arm $i_t^X \in [K]$.

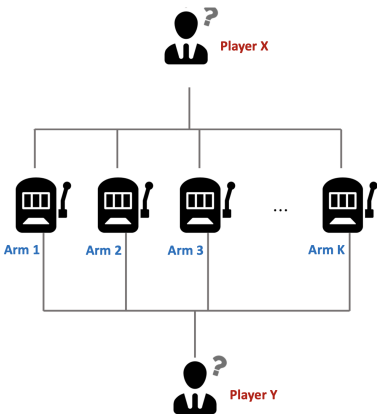
Multi-Player (Cooperative) Bandits



Consider $m > 1$ players with K stochastic bandit arms. Each time $1 \leq t \leq T$, each player X chooses an arm $i_t^X \in [K]$.

Players aim to maximize total reward (full cooperation). But they **cannot communicate** once the game starts.

Multi-Player (Cooperative) Bandits

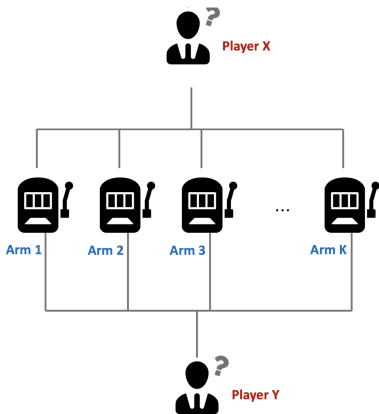


Consider $m > 1$ players with K stochastic bandit arms. Each time $1 \leq t \leq T$, each player X chooses an arm $i_t^X \in [K]$.

Players aim to maximize total reward (full cooperation). But they **cannot communicate** once the game starts.

Key challenge: colliding on the same action $i_t^X = i_t^Y$ at the same time yields zero reward.

Multi-Player (Cooperative) Bandits



Consider $m > 1$ players with K stochastic bandit arms. Each time $1 \leq t \leq T$, each player X chooses an arm $i_t^X \in [K]$.

Players aim to maximize total reward (full cooperation). But they **cannot communicate** once the game starts.

Key challenge: colliding on the same action $i_t^X = i_t^Y$ at the same time yields zero reward.

Proposed for wireless radio – learn good signal frequencies while avoiding interference.
[Lai-Jiang-Poor 08, Liu-Zhao 10, Anandkumar-Michael-Tang-Swami 11].

Partial Formulation

Fix $\mathbf{p} = (p_1, p_2, \dots, p_K) \in [0, 1]^K$. Generate TKm independent Bernoulli reward variables $\text{rew}_t^X(i)$ for $(t, i, X) \in [T] \times [K] \times [m]$:

$$\mathbb{P} \left[\text{rew}_t^X(i) = 1 \right] = p_i \quad \text{and} \quad \mathbb{P} \left[\text{rew}_t^X(i) = 0 \right] = 1 - p_i.$$

Partial Formulation

Fix $\mathbf{p} = (p_1, p_2, \dots, p_K) \in [0, 1]^K$. Generate TKm independent Bernoulli reward variables $\text{rew}_t^X(i)$ for $(t, i, X) \in [T] \times [K] \times [m]$:

$$\mathbb{P} \left[\text{rew}_t^X(i) = 1 \right] = p_i \quad \text{and} \quad \mathbb{P} \left[\text{rew}_t^X(i) = 0 \right] = 1 - p_i.$$

At time t , each player $(P_X)_{X \in [m]}$ picks arm i_t^X , and receives the reward:

$$\text{rew}_t(X) = \text{rew}_t^X(i_t^X) \cdot \mathbb{1}_{i_t^X \neq i_t^Y \quad \forall Y \neq X}.$$

Partial Formulation

Fix $\mathbf{p} = (p_1, p_2, \dots, p_K) \in [0, 1]^K$. Generate TKm independent Bernoulli reward variables $\text{rew}_t^X(i)$ for $(t, i, X) \in [T] \times [K] \times [m]$:

$$\mathbb{P} \left[\text{rew}_t^X(i) = 1 \right] = p_i \quad \text{and} \quad \mathbb{P} \left[\text{rew}_t^X(i) = 0 \right] = 1 - p_i.$$

At time t , each player $(P_X)_{X \in [m]}$ picks arm i_t^X , and receives the reward:

$$\text{rew}_t(\mathbf{X}) = \text{rew}_t^X(i_t^X) \cdot \mathbb{1}_{i_t^X \neq i_t^Y \quad \forall Y \neq X}.$$

$\mathbf{p}^* = \sum_{j=1}^m p_j^*$, the sum of the best m arms, is the regret benchmark:

$$R_T(\mathbf{p}) = \mathbb{E} \left[T\mathbf{p}^* - \left(\sum_{t=1}^T \sum_{X=1}^m \text{rew}_t(\mathbf{X}) \right) \right].$$

Partial Formulation

Fix $\mathbf{p} = (p_1, p_2, \dots, p_K) \in [0, 1]^K$. Generate TKm independent Bernoulli reward variables $\text{rew}_t^X(i)$ for $(t, i, X) \in [T] \times [K] \times [m]$:

$$\mathbb{P} \left[\text{rew}_t^X(i) = 1 \right] = p_i \quad \text{and} \quad \mathbb{P} \left[\text{rew}_t^X(i) = 0 \right] = 1 - p_i.$$

At time t , each player $(P_X)_{X \in [m]}$ picks arm i_t^X , and receives the reward:

$$\text{rew}_t(\mathbf{X}) = \text{rew}_t^X(i_t^X) \cdot \mathbb{1}_{i_t^X \neq i_t^Y \quad \forall Y \neq X}.$$

$\mathbf{p}^* = \sum_{j=1}^m p_j^*$, the sum of the best m arms, is the regret benchmark:

$$R_T(\mathbf{p}) = \mathbb{E} \left[T\mathbf{p}^* - \left(\sum_{t=1}^T \sum_{X=1}^m \text{rew}_t(\mathbf{X}) \right) \right].$$

Key quantities of interest: minimax and instance-dependent regret

$$R_T = \sup_{\mathbf{p} \in [0,1]^K} R_T(\mathbf{p}); \quad R_{T,\Delta} = \sup_{\mathbf{p} \in [0,1]^K: \mathbf{p}_m^* - \mathbf{p}_{m+1}^* \geq \Delta} R_T(\mathbf{p}).$$

Partial Formulation

Fix $\mathbf{p} = (p_1, p_2, \dots, p_K) \in [0, 1]^K$. Generate TKm independent Bernoulli reward variables $\text{rew}_t^X(i)$ for $(t, i, X) \in [T] \times [K] \times [m]$:

$$\mathbb{P} \left[\text{rew}_t^X(i) = 1 \right] = p_i \quad \text{and} \quad \mathbb{P} \left[\text{rew}_t^X(i) = 0 \right] = 1 - p_i.$$

At time t , each player $(P_X)_{X \in [m]}$ picks arm i_t^X , and receives the reward:

$$\text{rew}_t(\mathbf{X}) = \text{rew}_t^{X_t}(i_t^{X_t}) \cdot \mathbb{1}_{i_t^X \neq i_t^Y \quad \forall Y \neq X}.$$

$\mathbf{p}^* = \sum_{j=1}^m p_j^*$, the sum of the best m arms, is the regret benchmark:

$$R_T(\mathbf{p}) = \mathbb{E} \left[T\mathbf{p}^* - \left(\sum_{t=1}^T \sum_{X=1}^m \text{rew}_t(\mathbf{X}) \right) \right].$$

Key quantities of interest: minimax and instance-dependent regret

$$R_T = \sup_{\mathbf{p} \in [0,1]^K} R_T(\mathbf{p}); \quad R_{T,\Delta} = \sup_{\mathbf{p} \in [0,1]^K: \mathbf{p}_m^* - \mathbf{p}_{m+1}^* \geq \Delta} R_T(\mathbf{p}).$$

More specification needed! What information is observed about collisions?

Several Precise Formulations

Fix $\mathbf{p} = (p_1, p_2, \dots, p_K) \in [0, 1]^K$. Generate TKm independent Bernoulli reward variables $\text{rew}_t^X(i)$ for $(t, i, X) \in [T] \times [K] \times [m]$:

$$\mathbb{P} \left[\text{rew}_t^X(i) = 1 \right] = p_i \quad \text{and} \quad \mathbb{P} \left[\text{rew}_t^X(i) = 0 \right] = 1 - p_i.$$

At time t , each player $(P_X)_{X \in [m]}$ picks arm i_t^X , and receives the reward:

$$\text{rew}_t(X) = \text{rew}_t^X(i_t^X) \cdot \mathbb{1}_{i_t^X \neq i_t^Y \quad \forall Y \neq X}.$$

There are at least four natural feedback models when collisions occur.

Fix $\mathbf{p} = (p_1, p_2, \dots, p_K) \in [0, 1]^K$. Generate TKm independent Bernoulli reward variables $\text{rew}_t^X(i)$ for $(t, i, X) \in [T] \times [K] \times [m]$:

$$\mathbb{P} \left[\text{rew}_t^X(i) = 1 \right] = p_i \quad \text{and} \quad \mathbb{P} \left[\text{rew}_t^X(i) = 0 \right] = 1 - p_i.$$

At time t , each player $(P_X)_{X \in [m]}$ picks arm i_t^X , and receives the reward:

$$\text{rew}_t(X) = \text{rew}_t^X(i_t^X) \cdot \mathbb{1}_{i_t^X \neq i_t^Y \quad \forall Y \neq X}.$$

There are at least four natural feedback models when collisions occur.

- 1 **Strongly detectable:** the collision is explicitly announced.

Fix $\mathbf{p} = (p_1, p_2, \dots, p_K) \in [0, 1]^K$. Generate TKm independent Bernoulli reward variables $\text{rew}_t^X(i)$ for $(t, i, X) \in [T] \times [K] \times [m]$:

$$\mathbb{P} \left[\text{rew}_t^X(i) = 1 \right] = p_i \quad \text{and} \quad \mathbb{P} \left[\text{rew}_t^X(i) = 0 \right] = 1 - p_i.$$

At time t , each player $(P_X)_{X \in [m]}$ picks arm i_t^X , and receives the reward:

$$\text{rew}_t(X) = \text{rew}_t^X(i_t^X) \cdot \mathbb{1}_{i_t^X \neq i_t^Y \quad \forall Y \neq X}.$$

There are at least four natural feedback models when collisions occur.

- 1 **Strongly** detectable: the collision is explicitly announced.
- 2 **Weakly** detectable: observe realized reward 0.

Fix $\mathbf{p} = (p_1, p_2, \dots, p_K) \in [0, 1]^K$. Generate TKm independent Bernoulli reward variables $\text{rew}_t^X(i)$ for $(t, i, X) \in [T] \times [K] \times [m]$:

$$\mathbb{P} \left[\text{rew}_t^X(i) = 1 \right] = p_i \quad \text{and} \quad \mathbb{P} \left[\text{rew}_t^X(i) = 0 \right] = 1 - p_i.$$

At time t , each player $(P_X)_{X \in [m]}$ picks arm i_t^X , and receives the reward:

$$\text{rew}_t(X) = \text{rew}_t^X(i_t^X) \cdot \mathbb{1}_{i_t^X \neq i_t^Y \quad \forall Y \neq X}.$$

There are at least four natural feedback models when collisions occur.

- 1 **Strongly detectable:** the collision is explicitly announced.
- 2 **Weakly detectable:** observe realized reward 0.
- 3 **Undetectable:** observe “underlying” reward $\text{rew}_t^X(i_t^X)$.

Fix $\mathbf{p} = (p_1, p_2, \dots, p_K) \in [0, 1]^K$. Generate TKm independent Bernoulli reward variables $\text{rew}_t^X(i)$ for $(t, i, X) \in [T] \times [K] \times [m]$:

$$\mathbb{P} \left[\text{rew}_t^X(i) = 1 \right] = p_i \quad \text{and} \quad \mathbb{P} \left[\text{rew}_t^X(i) = 0 \right] = 1 - p_i.$$

At time t , each player $(P_X)_{X \in [m]}$ picks arm i_t^X , and receives the reward:

$$\text{rew}_t(X) = \text{rew}_t^X(i_t^X) \cdot \mathbb{1}_{i_t^X \neq i_t^Y \quad \forall Y \neq X}.$$

There are at least four natural feedback models when collisions occur.

- 1 **Strongly detectable**: the collision is explicitly announced.
- 2 **Weakly detectable**: observe realized reward 0.
- 3 **Undetectable**: observe “underlying” reward $\text{rew}_t^X(i_t^X)$.
- 4 **Adversarial**: observe a reward chosen by an adaptive adversary.

Regret for these Models

Strongly detectable: regret $\tilde{O}(\sqrt{T})$, even for non-stochastic. **Implicit communication.**
($\tilde{O}(\cdot)$ hides $\text{poly}(K, \log T)$ factors.) [Lugosi-Mehrabian 18, Bubeck-Li-Peres-S. 19]

Regret for these Models

Strongly detectable: regret $\tilde{O}(\sqrt{T})$, even for non-stochastic. **Implicit communication.** ($\tilde{O}(\cdot)$ hides $\text{poly}(K, \log T)$ factors.) [Lugosi-Mehrabian 18, Bubeck-Li-Peres-S. 19]

Weak detectability: regret $\tilde{O}(\sqrt{T})$ and $\tilde{O}(\frac{\log T}{\Delta})$. **Subtle implicit communication.** [Lugosi-Mehrabian 18, Huang-Combes-Trinh COLT 22, Pacchiano-Bartlett-Jordan 21]

Regret for these Models

Strongly detectable: regret $\tilde{O}(\sqrt{T})$, even for non-stochastic. **Implicit communication.** ($\tilde{O}(\cdot)$ hides $\text{poly}(K, \log T)$ factors.) [Lugosi-Mehrabian 18, Bubeck-Li-Peres-S. 19]

Weak detectability: regret $\tilde{O}(\sqrt{T})$ and $\tilde{O}(\frac{\log T}{\Delta})$. **Subtle implicit communication.** [Lugosi-Mehrabian 18, Huang-Combes-Trinh COLT 22, Pacchiano-Bartlett-Jordan 21]

What happens when communication is truly impossible? This is true already in (and a motivation for) the undetectable collision model.

Regret for these Models

Strongly detectable: regret $\tilde{O}(\sqrt{T})$, even for non-stochastic. **Implicit communication.** ($\tilde{O}(\cdot)$ hides $\text{poly}(K, \log T)$ factors.) [Lugosi-Mehrabian 18, Bubeck-Li-Peres-S. 19]

Weak detectability: regret $\tilde{O}(\sqrt{T})$ and $\tilde{O}(\frac{\log T}{\Delta})$. **Subtle implicit communication.** [Lugosi-Mehrabian 18, Huang-Combes-Trinh COLT 22, Pacchiano-Bartlett-Jordan 21]

What happens when communication is truly impossible? This is true already in (and a motivation for) the undetectable collision model.

Surprisingly, one can design **collision-free** algorithms attaining $\tilde{O}(\sqrt{T})$ regret (and using only public shared randomness). These automatically work in all feedback models.

Regret for these Models

Strongly detectable: regret $\tilde{O}(\sqrt{T})$, even for non-stochastic. **Implicit communication.** ($\tilde{O}(\cdot)$ hides $\text{poly}(K, \log T)$ factors.) [Lugosi-Mehrabian 18, Bubeck-Li-Peres-S. 19]

Weak detectability: regret $\tilde{O}(\sqrt{T})$ and $\tilde{O}(\frac{\log T}{\Delta})$. **Subtle implicit communication.** [Lugosi-Mehrabian 18, Huang-Combes-Trinh COLT 22, Pacchiano-Bartlett-Jordan 21]

What happens when communication is truly impossible? This is true already in (and a motivation for) the undetectable collision model.

Surprisingly, one can design **collision-free** algorithms attaining $\tilde{O}(\sqrt{T})$ regret (and using only public shared randomness). These automatically work in all feedback models.

Theorem (Bubeck-Budzinski 20, Bubeck-Budzinski-S. 21)

There exists an efficient, collision-free strategy with $\tilde{O}(\sqrt{T})$ regret. Precisely,

$$R_T = O\left(mK^{11/2}\sqrt{T \log T}\right),$$

$$\mathbb{P}(\text{there is ever a collision}) = O(T^{-2}).$$

Gap Dependent Regret Without Communication

Minimax regret barely changed. Turns out $R_{T,\Delta}$ changes a lot!

Minimax regret barely changed. Turns out $R_{T,\Delta}$ changes a lot!

Theorem (Liu-S. 22)

The Pareto optimal regret guarantees with no communication are:

$$R_{T,\Delta} \leq \tilde{O}\left(\frac{1}{\Delta_i \cdot \Delta_{i+1}}\right), \quad \Delta \in [\Delta_i, \Delta_{i+1}];$$
$$1 \geq \Delta_1 \geq \dots \geq \Delta_J \geq T^{-1/2}.$$

These are achievable with no collisions, hence in all feedback models.

Minimax regret barely changed. Turns out $R_{T,\Delta}$ changes a lot!

Theorem (Liu-S. 22)

The Pareto optimal regret guarantees with no communication are:

$$R_{T,\Delta} \leq \tilde{O}\left(\frac{1}{\Delta_i \cdot \Delta_{i+1}}\right), \quad \Delta \in [\Delta_i, \Delta_{i+1}];$$
$$1 \geq \Delta_1 \geq \dots \geq \Delta_J \geq T^{-1/2}.$$

These are achievable with no collisions, hence in all feedback models.

Extreme cases:

- $\Delta_J = 1$: minimax regret $R_T \leq \tilde{O}(\sqrt{T})$.
- $\Delta_j = 2^{-j}$: $R_{T,\Delta} \leq \tilde{O}(\Delta^{-2})$.

Minimax regret barely changed. Turns out $R_{T,\Delta}$ changes a lot!

Theorem (Liu-S. 22)

The Pareto optimal regret guarantees with no communication are:

$$R_{T,\Delta} \leq \tilde{O}\left(\frac{1}{\Delta_i \cdot \Delta_{i+1}}\right), \quad \Delta \in [\Delta_i, \Delta_{i+1}];$$
$$1 \geq \Delta_1 \geq \dots \geq \Delta_J \geq T^{-1/2}.$$

These are achievable with no collisions, hence in all feedback models.

Extreme cases:

- $\Delta_J = 1$: minimax regret $R_T \leq \tilde{O}(\sqrt{T})$.
- $\Delta_j = 2^{-j}$: $R_{T,\Delta} \leq \tilde{O}(\Delta^{-2})$.

Corollary: undetectable and adversarial models behave the same (up to $\text{poly}(K, \log T)$).

Corollary: if $\Delta \ll \Delta'$, no algorithm achieves $R_{T,\Delta} \leq \tilde{O}(1/\Delta)$ and $R_{T,\Delta'} \leq \tilde{O}(1/\Delta')$.

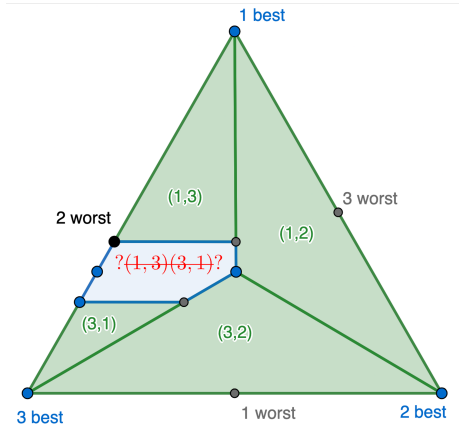
Geometric Viewpoint with 2 players and 3 actions

For illustration, work in the plane $P = \{p_1 + p_2 + p_3 = \text{constant}\}$ under full feedback. Undetectability means Player Y 's decisions do not influence Player X at all.

Geometric Viewpoint with 2 players and 3 actions

For illustration, work in the plane $P = \{p_1 + p_2 + p_3 = \text{constant}\}$ under full feedback.
Undetectability means Player Y 's decisions do not influence Player X at all.
Hence the protocol amounts to choosing for each $t \in [T]$ a function

$$(i_t^X, i_t^Y) : P \rightarrow \{1, 2, 3\}^2.$$



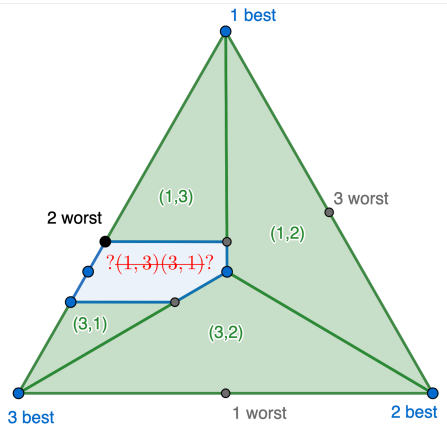
Geometric Viewpoint with 2 players and 3 actions

For illustration, work in the plane $P = \{p_1 + p_2 + p_3 = \text{constant}\}$ under full feedback.

Undetectability means Player Y 's decisions do not influence Player X at all.

Hence the protocol amounts to choosing for each $t \in [T]$ a function

$$(i_t^X, i_t^Y) : P \rightarrow \{1, 2, 3\}^2.$$



The estimates \hat{p}_t^X, \hat{p}_t^Y are within $\tilde{O}(t^{-1/2})$ of each other (by full feedback).

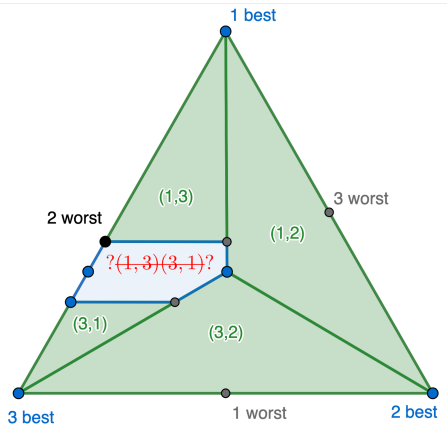
Geometric Viewpoint with 2 players and 3 actions

For illustration, work in the plane $P = \{p_1 + p_2 + p_3 = \text{constant}\}$ under full feedback.

Undetectability means Player Y 's decisions do not influence Player X at all.

Hence the protocol amounts to choosing for each $t \in [T]$ a function

$$(i_t^X, i_t^Y) : P \rightarrow \{1, 2, 3\}^2.$$

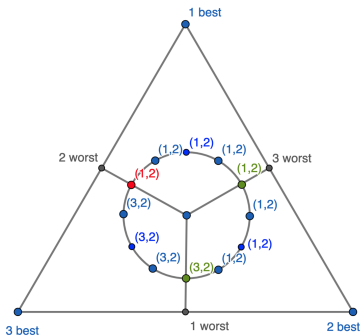


The estimates \hat{p}_t^X, \hat{p}_t^Y are within $\tilde{O}(t^{-1/2})$ of each other (by full feedback).

Difficulty: cannot always play the top 2 arms without colliding for some \mathbf{p} .

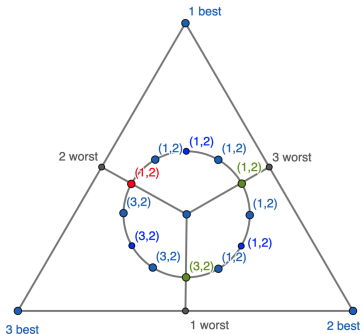
One-Step Regret Lower Bound with 2 players and 3 actions

How to turn this into a lower bound? Consider \sqrt{T} points equally spaced on a constant-size circle, labelled according to the time T strategy.



One-Step Regret Lower Bound with 2 players and 3 actions

How to turn this into a lower bound? Consider \sqrt{T} points equally spaced on a constant-size circle, labelled according to the time T strategy.

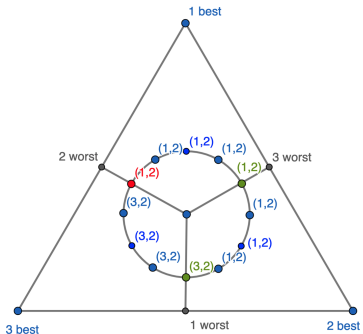


For any labelling, there is a **FAIL** incurring $\Omega(1)$ regret.
Either:

- 1 There is a collision, OR
- 2 The worst two actions are played.

One-Step Regret Lower Bound with 2 players and 3 actions

How to turn this into a lower bound? Consider \sqrt{T} points equally spaced on a constant-size circle, labelled according to the time T strategy.



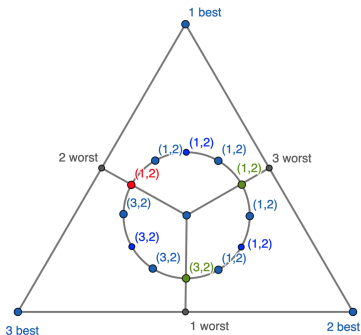
For any labelling, there is a **FAIL** incurring $\Omega(1)$ regret.
Either:

- 1 There is a collision, OR
- 2 The worst two actions are played.

The estimates $\hat{\mathbf{p}}_t^X, \hat{\mathbf{p}}_t^Y$ are basically adjacent points along this circle...

One-Step Regret Lower Bound with 2 players and 3 actions

How to turn this into a lower bound? Consider \sqrt{T} points equally spaced on a constant-size circle, labelled according to the time T strategy.



For any labelling, there is a **FAIL** incurring $\Omega(1)$ regret.
Either:

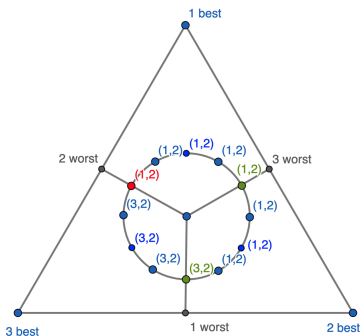
- 1 There is a collision, OR
- 2 The worst two actions are played.

The estimates $\hat{\mathbf{p}}_t^X, \hat{\mathbf{p}}_t^Y$ are basically adjacent points along this circle...

By dyadic pigeonhole, there **exists** a gap Δ_J with $\tilde{\Omega}(T)$ **FAILS** across $1 \leq t \leq T$.

One-Step Regret Lower Bound with 2 players and 3 actions

How to turn this into a lower bound? Consider \sqrt{T} points equally spaced on a constant-size circle, labelled according to the time T strategy.



For any labelling, there is a **FAIL** incurring $\Omega(1)$ regret. Either:

- 1 There is a collision, OR
- 2 The worst two actions are played.

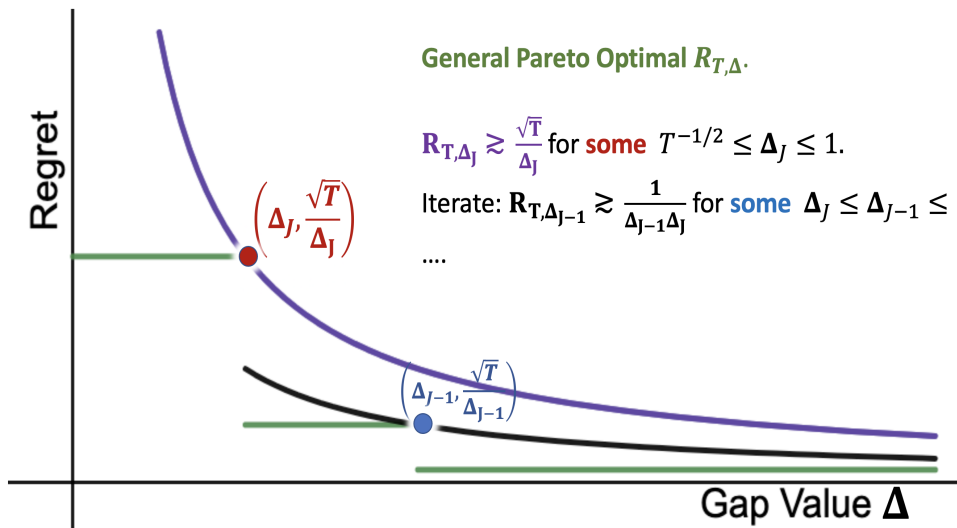
The estimates $\hat{\mathbf{p}}_t^X, \hat{\mathbf{p}}_t^Y$ are basically adjacent points along this circle...

By dyadic pigeonhole, there **exists** a gap Δ_J with $\tilde{\Omega}(T)$ **FAILs** across $1 \leq t \leq T$.

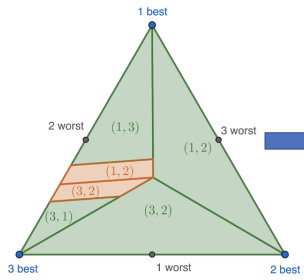
There are $\approx \Delta_J \sqrt{T}$ points on the circle with gap $\approx \Delta_J$ to absorb the **FAILs**. Hence

$$R_{T, \Delta_J} \gtrsim \frac{T}{\Delta_J \sqrt{T}} = \frac{\sqrt{T}}{\Delta_J}.$$

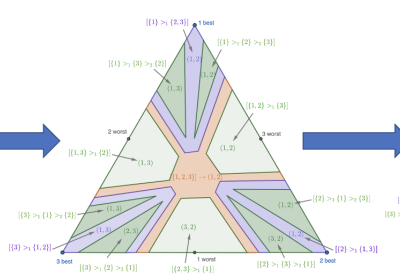
General Lower Bound: Set $T_J = \Delta_J^{-2}$ and Iterate



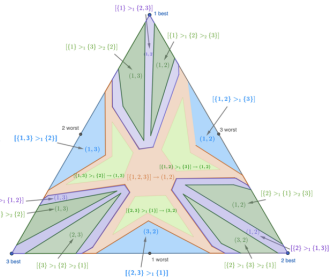
Collision-Free Algorithms At a Glance



[Bubeck-Budzinski 20]



[Bubeck-Budzinski-S 21]



[Liu-S 22]

- Previously: in multi-player stochastic bandits, $\tilde{O}(\sqrt{T})$ regret is possible with no collisions. Implicit communication enables $\tilde{O}(1/\Delta)$.
- This paper: without communication, Pareto optima include $\tilde{O}(\sqrt{T})$ and $\tilde{O}(1/\Delta^2)$. In particular, $\tilde{O}(1/\Delta)$ is only possible at a single scale Δ .

