# Crush Optimism with Pessimism:
# Structured Bandits Beyond Asymptotic Optimality

Kwang-Sung Jun (presenter)   and   Chicheng Zhang

THE UNIVERSITY OF ARIZONA

# Structured bandits

> e.g., linear $\mathcal{A} = \{a^1, \ldots, a^K \in \mathbb{R}^d\}$
> $\mathcal{F} = \{a \mapsto \theta^\top a : \theta \in \mathbb{R}^d\}$

- **Input**: Arm set $\mathcal{A}$, hypothesis class $\mathcal{F} \subset (\mathcal{A} \to \mathbb{R})$

"the set of possible **configurations** of the **mean rewards**"

- **Initialize**: The environment chooses $f^* \in \mathcal{F}$ (unknown to the learner)

  For $t = 1, \ldots, n$
  - Learner: chooses an arm $a_t \in \mathcal{A}$
  - Environment: generates the reward $r_t = f^*(a_t) + $ (zero-mean stochastic noise)
  - Learner: receives $r_t$

- **Goal**: Minimize the cumulative **regret**

$$\mathbb{E} \, \text{Reg}_n = \mathbb{E}\left[ n \cdot \left( \max_{a \in \mathcal{A}} f^*(a) \right) - \sum_{t=1}^n f^*(a_t) \right]$$

- Note: stochastic bandits with **realizability**

2

# Structured bandits

- **Why relevant?**
  Techniques developed here may extend to RL (e.g., ergodic RL [Ok+18])

- **Naive strategy: UCB**
  $\implies \dfrac{K}{\Delta}\log n$ regret bound (instance-dependent)
  - Scales with the number of arms $K$
  - Instead, the **complexity** of the hypothesis class $\mathcal{F}$ should appear.

- The **asymptotically optimal regret** is well-defined.
  - E.g., linear bandits : $c^* \cdot \log n$ for some well-defined $c^* \ll \dfrac{K}{\Delta}$.

| The goal of this paper |
|---|
| Achieve the **asymptotic optimality** with improved **finite-time** regret for any $\mathcal{F}$. |

# Asymptotic optimality

- **Optimism** in the face of uncertainty
  (e.g., UCB, Thompson sampling)

  $\implies$ optimal asymptotic / worst-case regret
  in $K$-**armed bandits**.

- Linear bandits: optimal worst-case rate = $d\sqrt{n}$

- Asymptotically optimal regret? $\implies$ **No!**
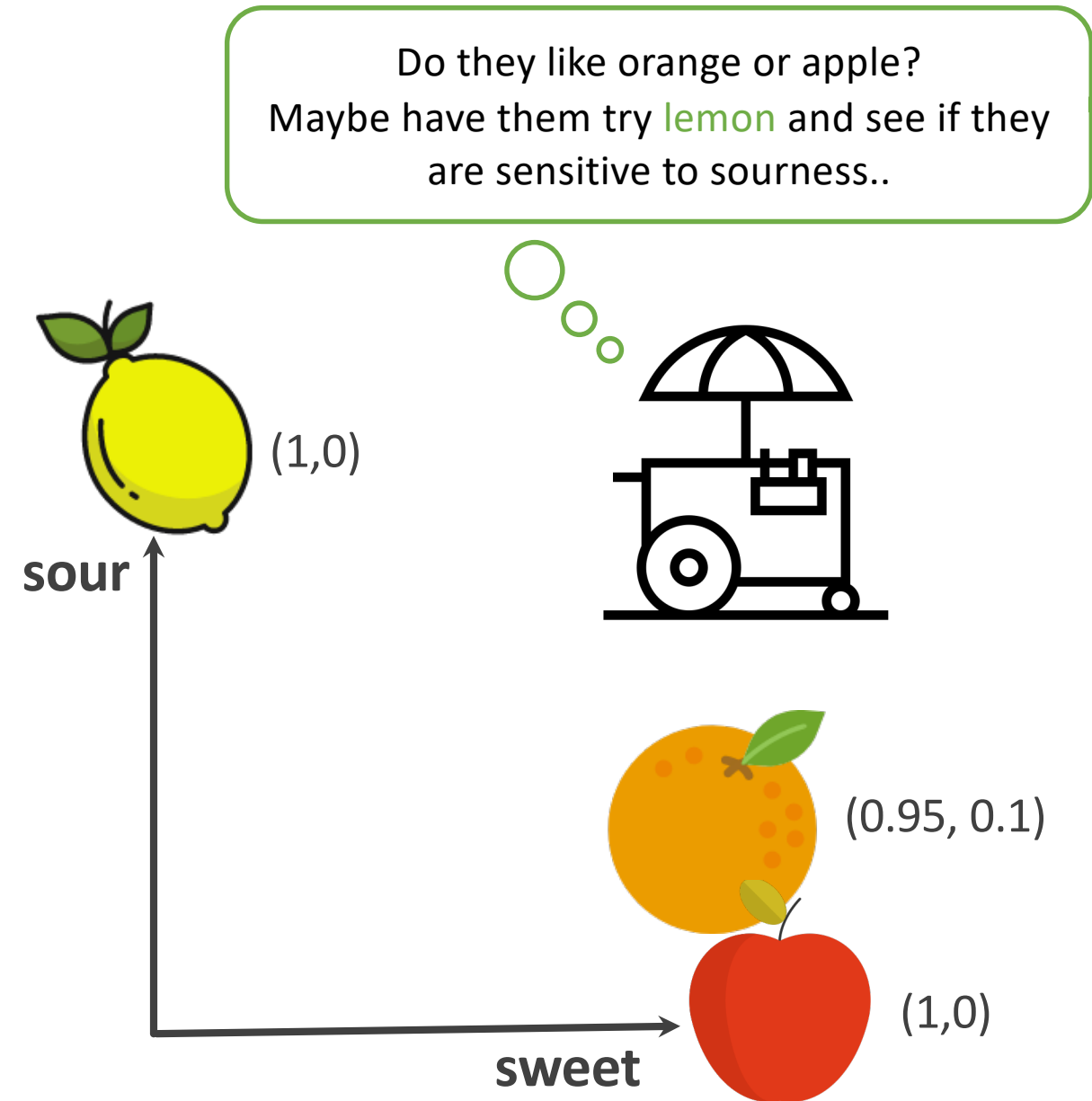
---

**The End of Optimism?**
**An Asymptotic Analysis of Finite-Armed Linear Bandits**

**Tor Lattimore**
Indiana University, Bloomington

**Csaba Szepesvári**
University of Alberta, Edmonton

(AISTATS'17)

Do they like orange or apple?
Maybe have them try lemon and see if they are sensitive to sourness..

**sour**

(1,0)

(0.95, 0.1)

(1,0)

**sweet**

# Asymptotic optimality: lower bound

- $\mathbb{E}\,\mathrm{Reg}_n \geq c(f^*) \cdot \log n$  (asymptotically)

$$\Delta_a = \left(\max_{b \in \mathcal{A}} f^*(b)\right) - f^*(a)$$

$$c(f^*) = \min_{\gamma_1,\dots,\gamma_K \geq 0} \sum_{a=1}^{K} \gamma_a \cdot \Delta_a$$

$$\text{s.\,t. } \forall g \in \mathcal{C}(f), \quad \sum_{a=1}^{K} \gamma_a \cdot \mathrm{KL}_\nu\big(f(a), g(a)\big) \geq 1$$

"competing" hypotheses

KL divergence with noise distribution $\nu$

- $\gamma^* = (\gamma_1^*, \dots, \gamma_K^*) \geq 0$ : the solution
- Suggests that we must pull arm $a$ like $\gamma_a^* \cdot \log n$ times.
- What if $c(f^*) = 0$? **Bounded regret!** (except for pathological ones)

# Existing asymptotically optimal algorithms

- Mostly uses forced exploration. [Lattimore+17,Combes+17,Hao+20]

  $\implies$ ensures **every arm**'s pull count is an **unbounded** function of $n$ such as $\frac{\log n}{1+\log \log n}$.

  $\implies \mathbb{E} \operatorname{Reg}_n \lesssim c(f^*) \cdot \log n + K \cdot \frac{\log n}{1+\log \log n}$

- Issues

  1. $K$ appears in the regret*   $\implies$   what if $K$ is exponentially large?

  2. **cannot** achieve **bounded** regret when $c(f^*) = 0$

- Parallel studies avoid forced exploration, but still depend on $K$. [Menard+20, Degenne+20]

*Dependence on $K$ can be avoided in special cases (e.g., linear).

# Contribution

| Research Question |
|---|
| Assume $\mathcal{F}$ is finite. Can we design an algorithm that<br><br>• enjoys the **asymptotic optimality**<br><br>• adapts to **bounded regret** whenever possible<br><br>• does not necessarily depend on $K$? |

Proposed algorithm:
**CRush Optimism with Pessimism (CROP)**

✓

✓

✓

- No forced exploration 😃
- The regret scales not with $K$ but with $K_\psi \leq K$ (defined in the paper).
- An interesting $\log \log n$ term in the regret*

# CROP, just the core part.

At time $t$,

- Maintain a confidence set $\mathcal{F}_t \subseteq \mathcal{F}$

- Do all $f \in \mathcal{F}_t$ agree on the best arm?
  - YES: pull that arm.
  - NO:

    - Compute the pessimism: $\overline{f}_t = \arg\min_{f \in \mathcal{F}_t} \max_{a \in \mathcal{A}} f(a)$

    - Compute $\gamma^* :=$ solution of the optimization problem $c\left(\overline{f}_t\right)$

    - (Tracking) Pull $a_t = \arg\min_{a \in \mathcal{A}} \dfrac{\text{pull\_count}(a)}{\gamma_a^*}$

> Cf. optimism: $\widetilde{f}_t = \arg\max_{f \in \mathcal{F}_t} \max_{a \in \mathcal{A}} f(a)$
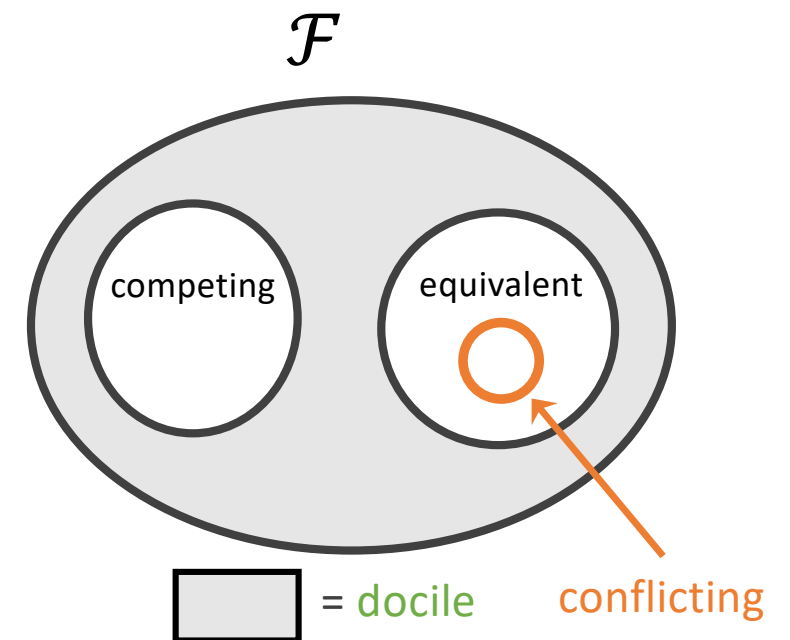
> Most existing approaches:
> Replace $\overline{f}_t$ with the empirical risk minimizer.
> $\implies$ requires forced sampling!

# Preview of the paper, no details.

$\mathcal{F}$

- Full version of CROP deals with
  - "docile" hypotheses. $\implies$ bounded terms in the regret
  - "conflicting" hypotheses $\implies \log \log n$ terms

- The **risk** of naively pursuing the asymptotic optimality
  - The oracle who plays according to $\gamma^*(f^*)$ may suffer a **linear regret** in the worst-case sense (finite time).

competing     equivalent

= docile     conflicting

- **It may not be the end of optimism**: we achieve **both** the worst-case and the asymptotic optimality by leveraging **optimism** (for a special $\mathcal{F}$ only).

Come to our poster after this session, find me in rocket chat, or email me/Chicheng for questions/discussions!