

Overview

Multi-armed bandit (MAB) problem is one of classic sequential decision-making problems with various real-world applications.

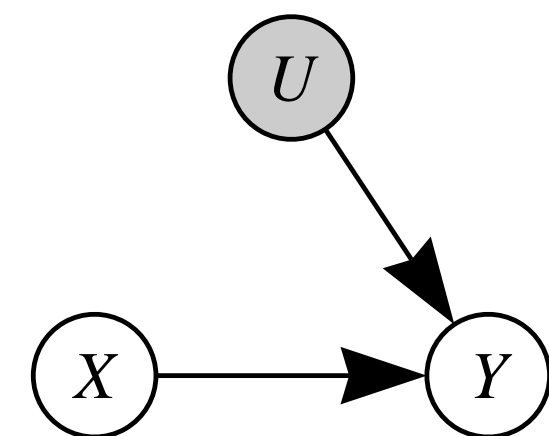
- **Arms:** There are arms \mathbf{A} in the bandit (i.e., slot machine); each arm associates with a reward distribution,
- **Play:** an agent plays the bandit by pulling an arm $A_{\mathbf{x}} \in \mathbf{A}$ each round,
- **Reward:** the bandit returns a reward $Y_{\mathbf{x}}$ drawn from the distribution,
- **Goal:** is to minimizing a cumulative regret (CR):

$$\text{Reg}_T = T\mu^* - \sum_{t=1}^T \mathbb{E}[Y_{A_t}] = \sum_{a=1}^K \Delta_a \mathbb{E}[T_a(T)]$$

where μ^* is the maximum expected reward

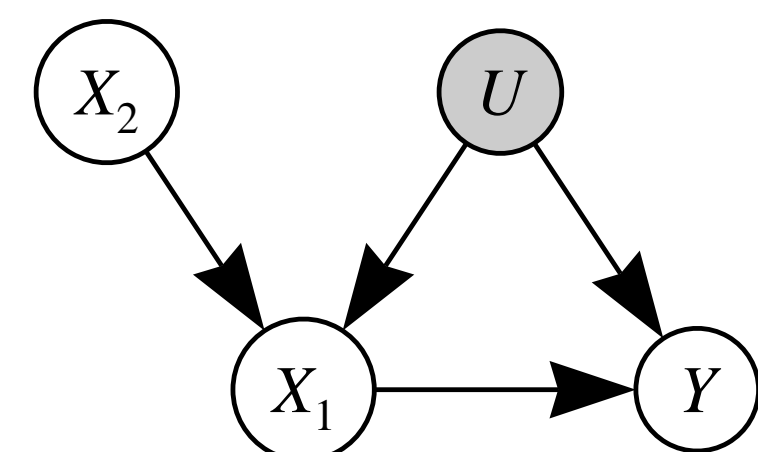
Multi-armed Bandit through Causal Lens

- pulling an arm = **intervening** a set of variables
- reward mechanism = **causal** mechanism

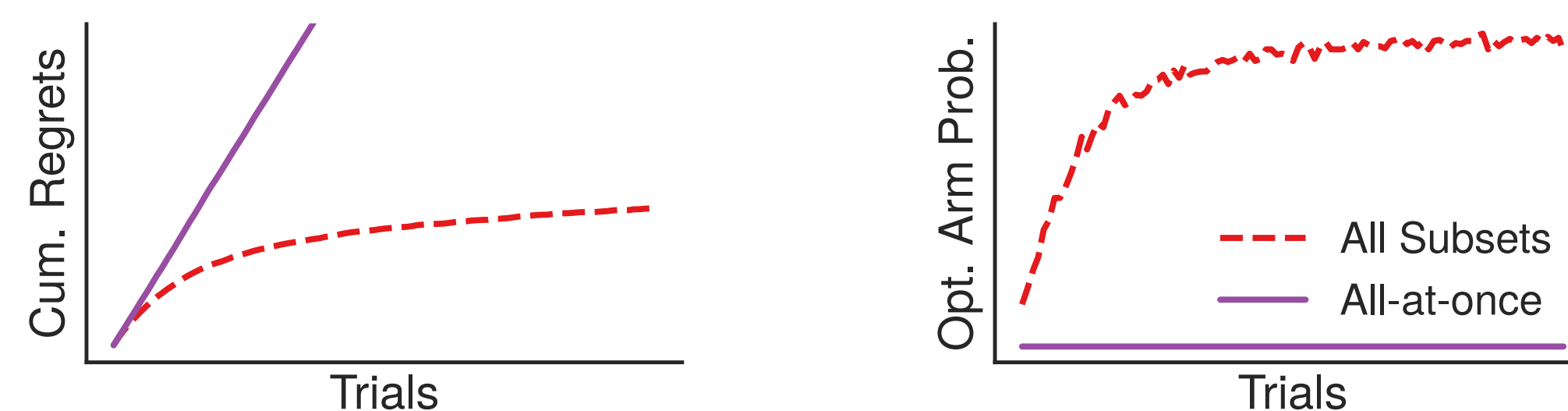


- Causally speaking, playing an arm A_x is setting X to x (called **do**), and observing Y drawn from $P(Y|do(X=x))$ where $P(y|do(x)) := \sum_u \mathbf{1}_{f(x,u),y} P(u)$.
- We are often oblivious to the existence of an underlying **causal mechanism**!

A Motivating Example: IV-MAB



- **Q:** How many **arms** are there? (You can control 2 binary variables, X_1 and X_2)
A: 9. setting values for $\{\emptyset, \{X_1\}, \{X_2\}, \{X_1, X_2\}\}$ (**all-subsets**). A naive combinatorial agent will intervene $\{X_1, X_2\}$ simultaneously (= 4 arms).
- **Q:** Why not *just* playing $\{X_1, X_2\}$ (**all-at-once**) altogether?
A: It might *miss* an optimal arm resulting:



There exists a case (i.e., parametrization) where intervening on X_2 is optimal, and intervening on $\{X_1, X_2\}$ simultaneously is always sub-optimal.
e.g., $X_1 = X_2 \oplus U$, $Y = X_1 \oplus U$. (when $X_2=1$, X_1 carries $\neg U$, and Y checks $X_1 \neq U$)

- **Q:** What are arms **worth** playing? (in *any* parametrizations)
A: intervening on either $\{X_2\}$ or $\{X_1\}$.

$$\therefore \max \mu_{X_2} \geq \max \mu_{\emptyset}, \quad \max \mu_{X_1} = \max \mu_{X_1, X_2}, \quad \max \mu_{X_2} < \max \mu_{X_1}$$

SCM-MAB — MAB built on the top of a causal framework

Structural Causal Model (SCM)

A Structural Causal Model (SCM) \mathcal{M} is a 4-tuple $\langle \mathbf{U}, \mathbf{V}, \mathbf{F}, P(\mathbf{U}) \rangle$:

- \mathbf{U} is a set of **unobserved** variables (**not modeled**);
- \mathbf{V} is a set of **observed** variables (**modeled**);
- \mathbf{F} is a set of **deterministic** functions for \mathbf{V} using \mathbf{U} and \mathbf{V} ;
- $P(\mathbf{U})$ is a joint distribution over the \mathbf{U} (**randomness**).

SCM-MAB

- SCM $\mathcal{M} = \langle \mathbf{U}, \mathbf{V}, \mathbf{F}, P(\mathbf{U}) \rangle$ and a reward variable $Y \in \mathbf{V}$, $\langle \mathcal{M}, Y \rangle$
- Arms \mathbf{A} correspond to *all* interventions $\{A_{\mathbf{x}} | \mathbf{x} \in D(\mathbf{X}), \mathbf{X} \subseteq \mathbf{V} \setminus \{Y\}\}$.
- Reward distribution $Y_{\mathbf{x}} \sim P(Y|do(\mathbf{X}=\mathbf{x}))$, expected reward, $\mu_{\mathbf{x}} := \mathbb{E}[Y|do(\mathbf{X}=\mathbf{x})]$.

We assume that a causal graph \mathcal{G} of \mathcal{M} is accessible not \mathbf{F} nor $P(\mathbf{U})$

Structural Properties in SCM-MAB — How arms are dependent

1. **Equivalence among Arms**

Two arms share the same reward distribution, e.g.,

$$\mu_{\mathbf{x}, \mathbf{z}} = \mu_{\mathbf{x}}$$

because intervening some of variables is redundant.

→ Test $P(y|do(\mathbf{x}, \mathbf{z})) = P(y|do(\mathbf{x}))$ through $Y \perp\!\!\!\perp \mathbf{Z} \mid \mathbf{X}$ in $\mathcal{G}_{\mathbf{X}\cup\mathbf{Z}}$ (*do-calculus*).

— Defn: **Minimal Intervention Set (MIS)**

Given that there are intervention sets with the same reward distribution, we would like to intervene a *minimal* set of variables yielding smaller # of arms.

2. **Partial-orderedness among Intervention Sets**

A set of variables \mathbf{X} may be preferred to the other set of variables \mathbf{Z} because their maximum achievable expected rewards can be ordered:

$$\mu_{\mathbf{x}^*} = \max_{\mathbf{x} \in D(\mathbf{X})} \mu_{\mathbf{x}} \geq \max_{\mathbf{z} \in D(\mathbf{Z})} \mu_{\mathbf{z}} = \mu_{\mathbf{z}^*}$$

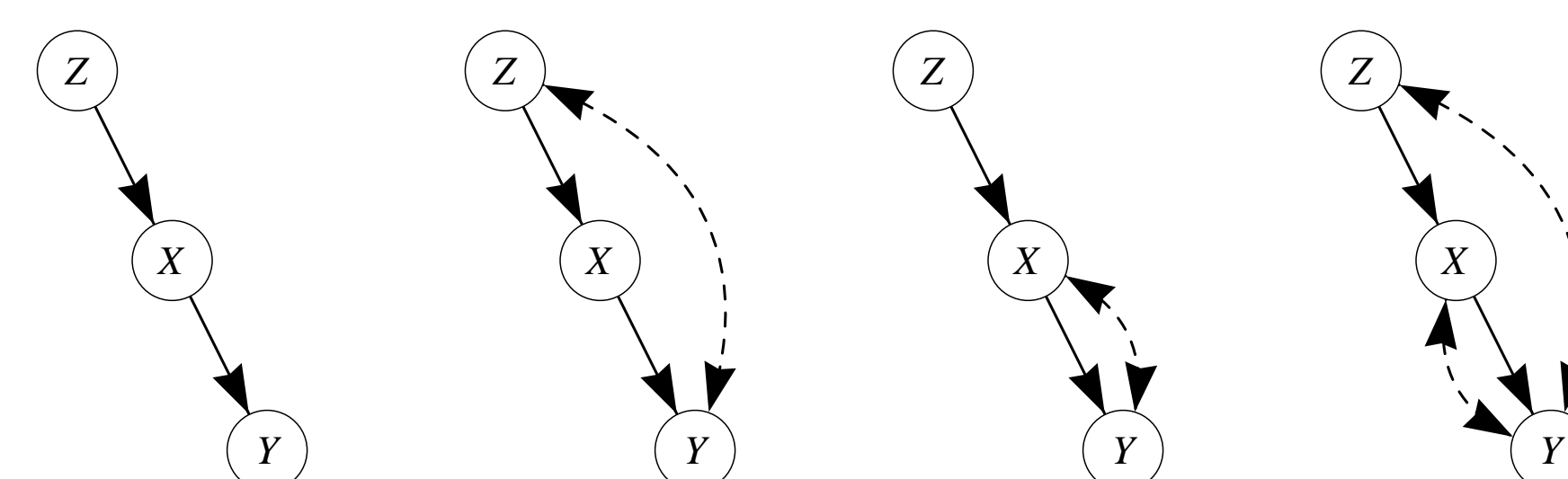
Pulling $do(\mathbf{z})$ delays the identification of optimal arms.

— Defn: **Possibly-Optimal Minimal Intervention Set (POMIS)**

Whether a set of variables, when intervened on, can yield **optimal expected reward** in *some* SCM \mathcal{M} conforming to the given causal graph \mathcal{G} .

Toy Examples for MISs and POMISs

(* a dashed bidirected edge = existence of an unobserved confounder)



Same **MISs** $\{\emptyset, \{X\}, \{Z\}\}$ since $do(x) = do(x, z)$ for $z \in D(\mathbf{Z})$.

POMIS are $\{\{X\}\}, \{\emptyset, \{X\}\}, \{\{Z\}, \{X\}\}, \{\emptyset, \{Z\}, \{X\}\}$

— We characterized a necessary and sufficient condition whether a set of variables is a **(PO)MIS**.

— We provide an algorithmic procedure to list all **(PO)MIS** given $\langle \mathcal{G}, Y \rangle$.

Empirical Evaluation

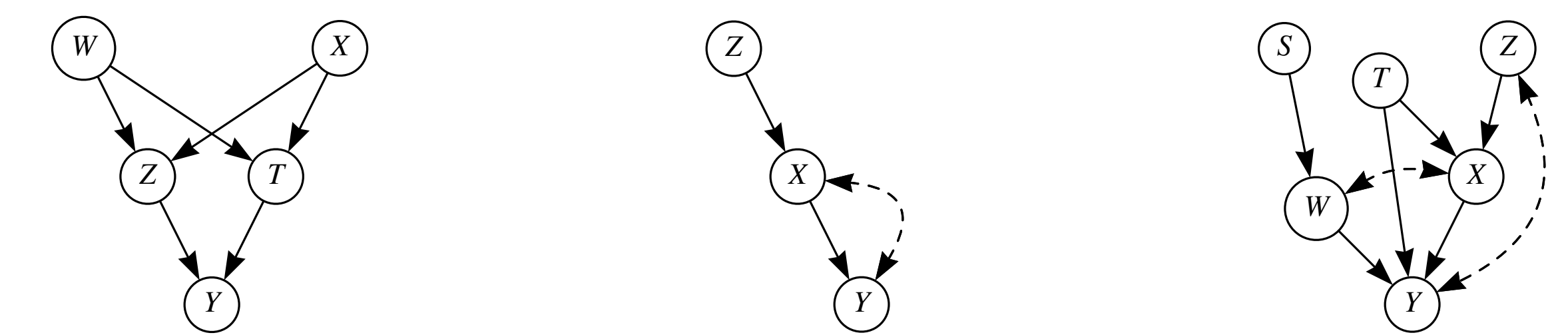
4 strategies \times 2 base MAB solvers \times 3 tasks; ($T = 10k$, 300 simulations)

Base MAB solvers: Thompson Sampling (TS) and kl-UCB

Strategies

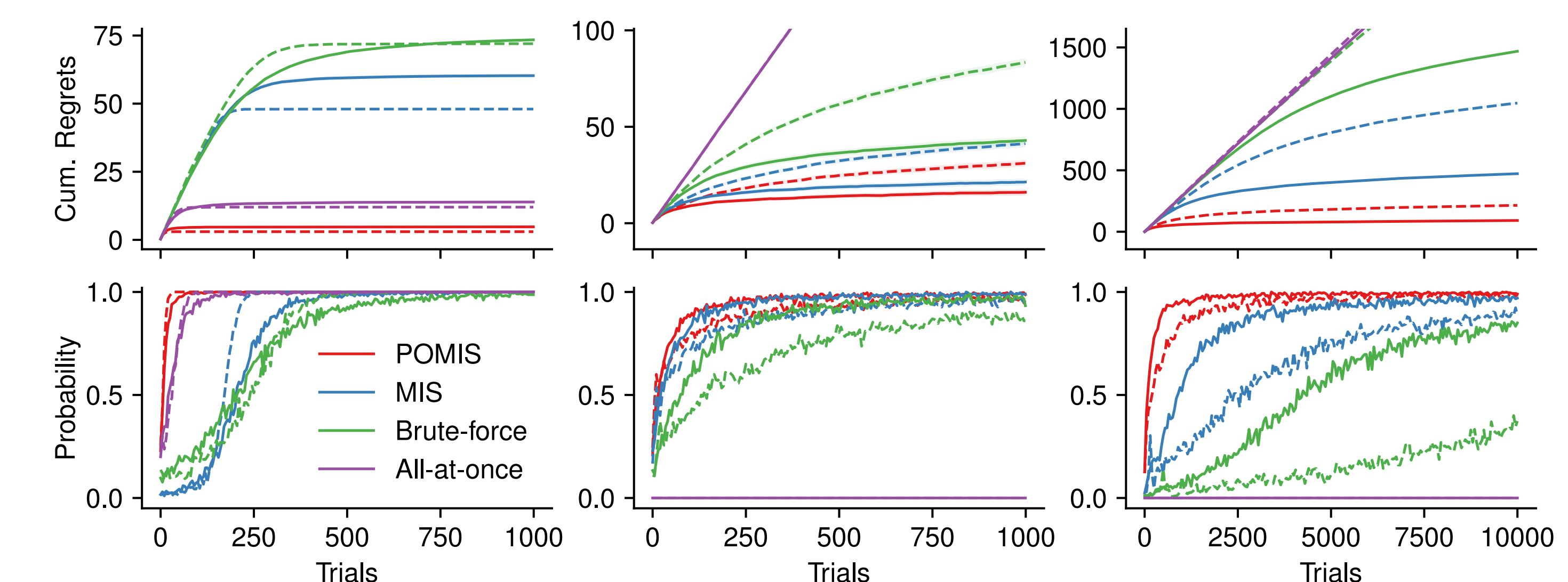
- **Brute-force:** all possible arms, $\{\mathbf{x} \in D(\mathbf{X}) \mid \mathbf{X} \subseteq \mathbf{V} \setminus \{Y\}\}$ (aka all-subsets)
- **All-at-once:** intervene all variables simultaneously, $D(\mathbf{V} \setminus \{Y\})$
- **MIS:** arms related to MISs
- **POMIS:** arms related to POMISs

Tasks



Results

(top) averaged cumulative regrets and (bottom) optimal arm probability
TS in solid lines, kl-UCB in dashed lines



CRs: **Brute-force** \geq **MIS** \geq **POMIS** (smaller the better)

If the number of arms for **All-at-once** is *smaller* than **POMIS**, then, it implies that **All-at-once** is missing possibly-optimal arms!

Conclusions

- introduced **SCM-MAB** = MAB + SCM = $\frac{\text{MAB}}{\text{SCM}}$
- characterized structural properties (equivalence, partial-orderedness) in SCM-MAB given a causal graph.
- studied conditions under which intervening a set of variables might lead to optimal! (POMIS)
- empirical results corroborate theoretical findings
- We have a *new* paper to be presented at **AAAI'2019**
- introduced **non-manipulability** constraints (not all variables are intervenable),
- characterized **MISs** / **POMISs** w/ the constraints,
- studied sophisticated relationships among POMISs arms

Code at <https://github.com/sanghack81/SCMMAB-NIPS2018>.

Papers at causalai.net.