# Combinatorial Multi-armed Bandits &

# $\alpha, \beta$ - Approximation

Lightning Day
September 15, 2023
Sharif University of Technology

Reza Pishkoo

## Presentation Overview

1. **Multi armed bandit (MAB)**

2. **Combinatorial multi armed bandit (CMAB)**

3. **General CMAB Framework**

4. **Main theorem**

5. **Comparing to classical MAB**

**1. Multi armed bandit (MAB)**

# Multi armed bandit

- m arms
- reward
- optimal arm
- regret
- tradeoff between exploration and exploitation

# real world application

- online advertising scenario
- website contains a set of web pages
- advertiser can select at most k web pages

# definition of CMAB

- super arms : any set of arms
- offline computation oracle computes the optimal super arm
- $(\alpha, \beta)$-approximation oracle
- $\alpha, \beta \leq 1$
- $\alpha$ : oracle outputs a super arm whose expected reward is at least $\alpha$ fraction of the optimal expected reward
- $\beta$ : success probability
- $(\alpha, \beta)$-approximation regret

## definition of CMAB framework

- m arms associated with a set of random variables $X_{i,t}$ for $1 \leq i \leq m$ and $t \geq 1$

- super arms : every set of arms

- $X_{i,t}$ : a random variable with bounded support on $[0,1]$

- $X_{i,t}$ : random outcome of the i-th arm in its t-th trial

- $\mu = (\mu_1, \mu_2, ..., \mu_m)$ : vector of expectations of all arms

- $T_{i,t}$ : number of times the outcome of arm i is revealed after the first t rounds

- $R_t(S)$ : non-negative random variable denoting the reward of super arm s

# Reward

- $\mathbb{E}\left[R_t(S)\right]$ : a function of only the set of arms S and the expectation vector $\mu$
- $r_\mu = \mathbb{E}\left[R_t(S)\right]$
- Monotonicity : if for all $i \in [m]$, $\mu_i \le \mu'_i$, we have $r_\mu(S) \le r_{\mu'}(S)$ for all S
- Bounded smoothness : there exists a stricly increasing function $f(.)$ such that, we have $|r_\mu(S) - r_{\mu'}(S)| \le f(\Lambda)$ if $\max_{i \in S} |\mu_i - \mu'_i| \le \Delta$

# CMAB algorithm

- Algorithm A selects $S_t^A$
- maximize $\mathbb{E}_{S,R}\left[\sum_{t=1}^n R_t(S_t^A)\right] = \mathbb{E}_S\left[\sum_{t=1}^n r_\mu(S_t^A)\right]$
- $opt_\mu = \max_S r_\mu(S)$
- $S_\mu^* = \max_S r_\mu(S)$
- $(\alpha, \beta)$-approximation oracle : $\mathbb{P}[r_\mu(S) \geq \alpha.opt_\mu] \geq \beta$
- $Reg_{\mu,\alpha,\beta}^A(n) = n.\alpha.\beta.opt_\mu - \mathbb{E}_S\left[\sum_{t=1}^n r_\mu(S_t^A)\right]$

# CMAB algorithm

- empirical mean :
  $\hat{\mu}_i = (\sum_{j=1}^{s} X_{i,j})/s$

- adjusted mean :
  $\bar{\mu}_i = \hat{\mu}_i + \sqrt{\frac{3ln(t)}{2T_i}}$

1: For each arm $i$, maintain: (1) variable $T_i$ as the total number of times arm $i$ is played so far; (2) variable $\hat{\mu}_i$ as the mean of all outcomes $X_{i,*}$'s of arm $i$ observed so far.
2: For each arm $i$, play an arbitrary super arm $S \in \mathcal{S}$ such that $i \in S$ and update variables $T_i$ and $\hat{\mu}_i$.
3: $t \leftarrow m$.
4: **while true do**
5:     $t \leftarrow t + 1$.
6:     For each arm $i$, set $\bar{\mu}_i = \hat{\mu}_i + \sqrt{\frac{3 \ln t}{2T_i}}$.
7:     $S = \text{Oracle}(\bar{\mu}_1, \bar{\mu}_2, \ldots, \bar{\mu}_m)$.
8:     Play $S$ and update all $T_i$'s and $\hat{\mu}_i$'s.
9: **end while**

**Algorithm 1:** CUCB with computation oracle

## Preliminaries

- Super arm S is bad if $r_\mu(S) < \alpha.opt_\mu$

- $S_B = \{S|r_\mu(S) < \alpha.opt_\mu\}$

- $\Delta_{min}^i = \alpha.opt_\mu - \max\{r_\mu(S)|S \in S_B, i \in S\}$

- $\Delta_{max}^i = \alpha.opt_\mu - \min\{r_\mu(S)|S \in S_B, i \in S\}$

- $\Delta_{max} = \max_{i \in [m]} \Delta_{max}^i$

- $\Delta_{min} = \min_{i \in [m]} \Delta_{min}^i$

# Theorem

## Theorem

*The $(\alpha, \beta)$-approximation regret of the CUCB algorithm in n rounds using an $(\alpha, \beta)$-approximation oracle is at most*

$$\sum_{i \in [m], \Delta_{min}^i > 0} \left( \frac{6ln(n)\Delta_{min}^i}{\left(f^{-1}(\Delta_{min}^i)\right)^2} + \int_{\Delta_{min}^i}^{\Delta_{max}^i} \frac{6ln(n)}{\left(f^{-1}(x)\right)^2} dx \right) + \left( \frac{\pi^2}{3} + 1 \right) . m . \Delta_{max}$$

*simplified form :*

$$\left( \frac{6ln(n)}{\left(f^{-1}(\Delta_{min})\right)^2} + \frac{\pi^2}{3} + 1 \right) . m . \Delta_{max}$$

# Comparing to classical MAB

- every super arm is a simple arm

- $f(x) = x$

- $\alpha = \beta = 1$

- $\Delta^i_{max} = \Delta^i_{min}$

- $\Delta^i = \max_{j \in [m]} \mu_j - \mu_i$

  the regret bound of the classical MAB

$$\sum_{i \in [m], \Delta^i > 0} \frac{6 ln(n)}{\Delta^i} + \left( \frac{\pi^2}{3} + 1 \right) . m . \Delta_{max}$$

- the coefficient of the regret bound in (Auer et al., 2002) : $\sum_{i \in [m], \Delta^i > 0} \frac{8}{\Delta^i}$