Background and motivation

Setting and algorithm 00000

Lower bound

Upper bound

Conclusion and future work

The Hardness Analysis of Thompson Sampling for Combinatorial Semi-bandits with Greedy Oracle

Fang Kong¹, Yueran Yang¹, Wei Chen², Shuai Li¹

¹John Hopcroft Center for Computer Science, Shanghai Jiao Tong University

²Microsoft Research Asia

NeurIPS 2021





イロト 不得 トイヨト イヨト

-

Background and motivation •00000000000 Setting and algorithm

Lower bound 0000 Upper bound

Conclusion and future work

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Stochastic multi-armed bandit



- T-round learning game between the player and m arms
- each arm i ∈ [m] is associated with an unknown fixed reward distribution D_i with unknown mean μ_i

Background and	motivation 00	Settin 0000	g and a O	lgorithr	n	Lower 0000	bound		Uppei 000	boui	nd	Cone	lusion	and futu	re work	
	S	Stoc	chas	stic	mu	Iti-	arn	ned	ba	nc	lit					
	Time Left arm	1 \$1	2 \$0	3	4	5 \$1	6 \$1	7 \$0	8	9	10	11	12			

Right arm \$1 \$0

• for round t = 1, 2, ...

- the player selects an arm $A_t \in [m]$
- then receives a reward $X_{t,A_t} \sim D_{A_t}$
- The goal of the player is to maximize the cumulative expected reward, or equivalent to minimizing the cumulative expected regret

$$\mathbb{E}\left[T\cdot\mu^*-\sum_{t=1}^T X_{t,A_t}\right] \tag{1}$$

Background and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work
00000000000	00000	0000	000	000

Stochastic multi-armed bandit

To achieve this long-horizon goal:

Exploration V.S. Exploitation

- exploration: try arms that have not been observed enough times
- exploitation: focus on the arm with the best observed performance so far

How to balance exploration and exploitation is the key of algorithms

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●



upper confidence bound (UCB)

- construct confidence sets for unknown expected rewards
- select arms according to their highest upper confidence bounds
- regret upper bound of order $O\left(\frac{\log T}{\Delta}\right)^{-1}$



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

$${}^{1}\Delta = \min_{i \in [m]} \left\{ \mu^{*} - \mu_{i} \right\}$$



Thompson sampling (TS)

- maintain an iteratively updated posterior distribution for unknown expected rewards
- select arms according to probabilities of being the best one
- was introduced in 1933², but has not been theoretically proven until recent years³
- regret upper bound of order $O\left(\frac{\log T}{\Delta}\right)$



 2 William R Thompson. "On the likelihood that one unknown probability exceeds another in view of the evidence of two samples". In: *Biometrika* 25.3/4 (1933), pp. 285–294.

Combinatorial multi-armed bandit (CMAB)

Advertisement placement

- The agent needs to select several webpages (left) to place the advertisement, thus the click-through rate is maximized when faced with several users (right).
- can be formulated by probabilistic maximum coverage (PMC)



Background and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work
000000000000	00000	0000	000	000

Combinatorial multi-armed bandit (CMAB)

- for round *t* = 1, 2, . . .:
- the player selects a combination of arms: S_t satisfying $S_t \in C$
- the environment samples the outcome $X_{t,i}$ for each arm i
- the player observes the feedback $Q_t = \{(i, X_{t,i}) : i \in S_t\}$ and receives the reward $R(S_t, X_t)$
- The goal of the player is to maximize the cumulative expected reward

$$\mathbb{E}_{S_t,X_t}\left[\sum_{t=1}^T R(S_t,X_t)\right]$$

Background and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work

Assumption

The expected reward of an action S only depends on the S and the mean vector μ . That is to say, there exists a function r such that $\mathbb{E}[R_t] = \mathbb{E}_{X_t \sim D}[R(S_t, X_t)] = r(S_t, \mu).$

Assumption

(Lipschitz continuity) There exists a constant B such that for any action S and mean vectors μ, μ' , the reward of S under μ and μ' satisfies

$$\left| \mathsf{r}(\mathsf{S},\mu) - \mathsf{r}(\mathsf{S},\mu') \right| \leq B \sum_{i\in\mathsf{S}} \left| \mu_i - \mu'_i \right| \,.$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Combinatorial multi-armed bandit (CMAB)

CMAB algorithms usually rely on an oracle to solve the corresponding offline problem.

- when the expected rewards are known, the problem of finding the optimal solution $\operatorname{argmax}_{S\in\mathcal{C}} r(S,\mu)$ is called *offline* problem
- many offline problems are NP-hard (such as maximum coverage, influence maximization) and only approximate algorithms (oracles) are available
- an Oracle is (α, β)-approximate if
 r(Oracle(μ'), μ') ≥ α max_{S∈C} r(S, μ') with probability larger
 than β for any given μ'

Background and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work
00000000000000	00000	0000	000	000

Combinatorial multi-armed bandit (CMAB)

With an (α, β) -approximate oracle, the goal is equivalent to minimize the (α, β) -approximate regret

$$\mathbb{E}\left[T \cdot \alpha\beta \cdot r(S^*, \mu) - \sum_{t=1}^T r(S_t, \mu)\right]$$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Background and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work
000000000000000	00000	0000	000	000

Related work: two popular algorithms

upper confidence bound - type (UCB)

- construct confidence sets for unknown expected reward of each arm
- an oracle helps to choose the action according to arms' highest upper confidence bounds
- require the reward function to satisfy the monotonicity on μ

Algorithm type	allow approximate oracle	regret guarantee ⁴
UCB-type	yes, $orall lpha, eta \in (0,1]$	$O\left(\frac{\log T}{\Delta}\right)^5$

 ${}^{4}\Delta = \min_{S \in \mathcal{C}: S \notin \alpha \text{ OPT }} \left\{ \alpha \cdot r(S^{*}, \mu) - r(S, \mu) \right\}$

 5 Qinshi Wang and Wei Chen. "Improving regret bounds for combinatorial semi-bandits with probabilistically triggered arms and its applications". In: Advances in Neural Information Processing Systems. 2017, pp. 1461–1121.

Background and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work
000000000000	00000	0000	000	000

Related work: two popular algorithms

Thompson sampling - type (TS)

- maintain an iteratively updated posterior distribution for unknown expected reward of each arm
- an oracle helps to choose the action according to sampled parameters (θ_i)_{i∈[m]} from posterior distributions
- relax the monotonicity requirment of UCB-type algorithms
- easier implementation and better empirical performance

Algorithm type	allow approximate oracle	regret guarantee ⁶
TS-type	no, $lpha=eta=1$	$O\left(\frac{\log T}{\Delta}\right)^7$

 ${}^{6}\Delta = \min_{S \in \mathcal{C}: S \notin OPT} \left\{ r(S^{*}, \mu) - r(S, \mu) \right\}$

⁷Pierre Perrault et al. "Statistical Efficiency of Thompson Sampling for Combinatorial Semi-Bandits". In: Advances in Neural Information Processing Systems. 2020, Siwei Wang and Wei Chen. "Thompson sampling for combinatorial semi-bandits". In: International Conference on Machine Learning. cPMLR 2018 (pp. 5114-5122.

Background and motivation 00000000000●	Setting and algorithm	Lower bound	Upper bound 000	Conclusion and future work
	Mot	ivation		

An example [Theorem 2, $Work^8$] shows the failure of TS to learn with an approximation oracle.

But this oracle is uncommon and is designed only for a specific problem instance.

Can the convergence analysis of TS extend beyond the exact oracle in CMAB?

⁸Siwei Wang and Wei Chen. "Thompson sampling for combinatorial semi-bandits". In: International ▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00 Conference on Machine Learning. PMLR. 2018, pp. 5114-5122.

Background and motivation	Setting and algorithm ●0000	Lower bound	Upper bound	Conclusion and future work
	Se	etting		

We consider TS with common Greedy oracle, which can provide approximate solutions for many offline problems

- *m* base arms, denoted by $[m] = \{1, 2, \dots, m\}$
- due to the problem structure, arms are further divided into n units, a unit of arms must be selected together. Denote the unit set as U
- for round $t = 1, 2, \ldots$
- the player selects action $S_t \in \mathcal{S} = \{S \subseteq \mathcal{U} : |S| = K\}$

- the environment samples the outcome $X_{t,i}$, $\forall i \in [m]$
- the player observes the feedback
 - $Q_t = \{(i, X_{t,i}) : i \in s \text{ for some } s \in S_t\}$
- the player receives the reward $R(S_t, X_t)$

Background and motivation	Setting and algorithm ○●○○○	Lower bound	Upper bound 000	Conclusion and future work

Greedy oracle

Algorithm 1 Greedy algorithm

- Input: base arm set [m] and mean vector μ = (μ_i)_{i∈[m]}, unit set U, action size K
- 2: Initialize: $S_{g,0} = \emptyset$
- 3: for $k = 1, 2, \cdots, K$ do
- 4: $s_k = \operatorname{argmax}_{s \in \mathcal{U} \setminus S_{g,k-1}} r(S_{g,k-1} \cup \{s\}, \mu)$
- 5: $S_{g,k} = S_{g,k-1} \cup \{s_k\}$
- 6: end for
- 7: Output: $S_g = S_{g,K}$

The above framework with Greedy oracle can cover many CMAB problems including probabilistic maximum coverage (PMC).

For simplicity, we assume S_g is unique under μ .

ground and motivation	Setting and algorithm ○○●○○	Lower bound	Upper bound	Conclusion and future work

Combinatorial Thompson sampling (CTS)

Algorithm 2 CTS algorithm with Greedy oracle

- 1: Input: base arm set [m], unit set \mathcal{U} , action size K
- 2: Initialize: $\forall i \in [m], G_i$ is the prior distribution
- 3: for $t = 1, 2, \cdots$ do
- 4: $\forall i \in [m]$: Sample $\theta_{t,i} \sim G_i$.
- 5: Select action $S_t = Greedy([m], \theta_t, U, K)$ and receive the observation Q_t

A D N A 目 N A E N A E N A B N A C N

- 6: //update
- 7: for $(i, X_{t,i}) \in Q_t$ do
- 8: Update the posterior distribution G_i
- 9: end for

10: end for

kground and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work

Objective: α -approximate regret?

under UCB:

$$egin{aligned} lpha \cdot r(S^*,\mu) - r(S_t,\mu) &\leq lpha \cdot r(S^*,U) - r(S_t,\mu) \ &\leq r(S_t,U) - r(S_t,\mu) \ &\leq \sum_{i\in \mathcal{S}_t} |U_i - \mu_i| \,. \end{aligned}$$

under TS:

$$\begin{aligned} &\alpha \cdot r(S^*, \mu) - r(S_t, \mu) \\ &= \alpha \cdot r(S^*, \mu) - \alpha \cdot r(S^*, \theta) + \alpha \cdot r(S^*, \theta) - r(S_t, \mu) \\ &\leq \alpha \cdot (r(S^*, \mu) - r(S^*, \theta)) + r(S_t, \theta) - r(S_t, \mu) \\ &\leq \alpha \sum_{i \in S^*} |\theta_i - \mu_i| + \sum_{i \in S_t} |\theta_i - \mu_i| \,. \end{aligned}$$

It is first brought up in analyzing UCB-based algorithms and may not well fit TS-based algorithms.

Background and motivation	Setting and algorithm 0000●	Lower bound	Upper bound 000	Conclusion and future work
	Ob	jective		

A different objective to minimize the greedy regret

$$R_g(T) = \mathbb{E}\left[\sum_{t=1}^T \max\left\{r(S_g, \mu) - r(S_t, \mu), 0\right\}\right].$$
 (2)

Remark

When Greedy can provide α -approximate solutions, the upper bound for the greedy regret also implies the upper bound for the α -approximate regret since $r(S_g, \mu) \ge \alpha \cdot r(S^*, \mu)$.

Background and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work
000000000000	00000	0000	000	000

Lower bound: define gaps to quantify the hardness

For any unit $s \in U$ and index $k \in [K]$ such that $s \notin S_{g,k}$, define the marginal reward gap

$$\Delta_{\boldsymbol{s},\boldsymbol{k}} = \boldsymbol{r}(S_{\boldsymbol{g},\boldsymbol{k}},\mu) - \boldsymbol{r}(S_{\boldsymbol{g},\boldsymbol{k}-1} \cup \{\boldsymbol{s}\},\mu)$$

as the reward difference between $S_{g,k}$ and $S_{g,k-1} \cup \{s\}$. According to the Greedy algorithm, we have $\Delta_{s,k} > 0$ for any k such that $s \notin S_{g,k}$. And for any action $S \in S$, define $\Delta_S = \max \{r(S_g, \mu) - r(S, \mu), 0\}$ as the reward difference from the Greedy's solution S_g . Let

$$\Delta_{s}^{\min} = \min_{S \in \mathcal{S}: s \in S} \Delta_{S} , \quad \Delta_{s}^{\max} = \max_{S \in \mathcal{S}: s \in S} \Delta_{S}$$

be the minimum and maximum reward gap of actions containing unit *s*, respectively.

- ロ ト - 4 回 ト - 4 □

Background and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work
000000000000	00000	0000	000	000

Lower Bound: problem instance



Figure: The underlying graph of the PMC instance used to derive the hardness analysis and the corresponding rewards of actions.

$$S_g = \{u_2, u_1\}$$
 with $s_{g,1} = u_2$, $s_{g,2} = u_1$, and $r(S_g, \mu) = 0.892$, while the optimal action is $\{u_1, u_4\}$

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Background and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work
000000000000	00000	0000	000	000

Lower bound: intuition and result

Using the CTS algorithm with Gaussian priors and Greedy oracle to solve the CMAB problem shown in Figure 1, when T is sufficiently large, we have

$$\mathbb{E}\left[N_{T+1,s}\right] = \Omega\left(\frac{\log T}{\Delta_{s,1}^2}\right),\qquad(3)$$

for any $s \neq s_{g,1} = u_2$, where $N_{T+1,s} = \sum_{t=1}^{T} \mathbb{1}\{s \in S_t\}$ is the number of rounds when s is contained in the selected action set S_t .

Background and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work
000000000000	00000	0000	000	000

Lower bound: intuition and result

$$\Delta_{u_3,1} = \Delta, \quad \Delta_{u_3,2} = 0.012 + 0.7\Delta, \quad \Delta_{u_3}^{\mathsf{min}} = 0.52\Delta + 0.008 \,.$$

Theorem

(Lower bound) When T is sufficiently large, the cumulative greedy regret satisfies

$$R_g(T) = \Omega\left(\frac{\log T}{\Delta_{u_3,1}^2} \cdot \Delta_{u_3}^{\min}\right) = \Omega\left(\frac{\log T}{\Delta^2}\right).$$
(4)

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Background and motivation	Setting and algorithm	Lower bound	Upper bound ●00	Conclusion and future work

Upper bound

Theorem

(Upper bound) The cumulative greedy regret of CTS with Gaussian and Beta priors can be upper bounded by

$$R_g(T) = O\left(\sum_{s \neq s_{g,1}} \max_{k:s \notin S_{g,k}} \frac{B^2 |s|^2 \Delta_s^{\max} \log T}{\Delta_{s,k}^2}\right), \quad (5)$$

where B is the coefficient of the Lipschitz continuity condition, $|\cup S_g|$ is the number of base arms that belong to the units contained in S_g .



Upper bound: Comparisons with TS for MAB

- K = 1, B = 1 and |s| = 1 for any s
- Greedy can provide exact optimal solutions thus $R_g(T) = R(T) = \mathbb{E}\left[\sum_{t=1}^{T} (\mu^* - \mu_{A_t})\right]$

•
$$\Delta_s^{\max} = \Delta_{s,1}$$
 for any s

Our greedy regret upper bound is

$$R_g(T) = O\left(\sum_{s \neq s_{g,1}} \frac{\log T}{\Delta_{s,1}}\right),$$

which recovers the main order of the regret of TS for MAB⁹.

⁹Shipra Agrawal and Navin Goyal. "Further optimal regret bounds for thompson sampling". In: Proceedings of the 16th International Conference on Artificial Intelligence and Statistics. 2013, pp>99–107. < ≥ > < ≥ > < ≥ < <>>

Background and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work
000000000000	00000	0000	000	000

Upper bound: proof sketch

Above all, each unit s needs to be explored

$$O\left(\max_{k:s\notin S_{g,k}}\frac{\log T}{\Delta_{s,k}^2}\right)$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Background and motivation	Setting and algorithm	Lower bound	Upper bound 000	Conclusion and future work ●○○
	Con	clusion		

- The first theoretical results for TS-type algorithm to solve CMAB problems with (approximate) greedy oracle.
- Our results break the misconception that CTS cannot be used with approximation oracles with the failure example in work¹⁰.

¹⁰Siwei Wang and Wei Chen. "Thompson sampling for combinatorial semi-bandits". In: *International* Conference on Machine Learning. PMLR. 2018, pp. 5114–5122. (□ > (□ > (□) + (0) + (0)

Background and motivation	Setting and algorithm	Lower bound	Upper bound	Conclusion and future work ○●○
	Futu	ire work		

An interesting future direction is to extend the current CMAB framework to the case with probabilistically triggered arms (CMAB-T).

- the CMAB-T framework models more applications
- TS samples candidate parameters to escape the computation of complicated optimization problems which may be faced in UCB¹¹
- new proof techniques are required

Background and motivation

Setting and algorithm 00000

Lower bound 0000 Upper bound

Conclusion and future work $_{\bigcirc \bigcirc \bigcirc}$

(日) (個) (E) (E) (E)

Thanks! Questions?