# The Symmetry between Arms and Knapsacks: A Primal-Dual Approach for Bandits with Knapsacks

Xiaocheng Li[1], Chunlin Sun[2], Yinyu Ye[3]

[1]Imperial College Business School
[2]Institute for Computational and Mathematical Engineering, Stanford University
[3]Department of Management Science and Engineering, Stanford University

## 1. Multi-Armed Bandit with Knapsack



- **Bandit:** Consider the bandit machine has $m$ arms with fixed but unknown mean reward $\mu_i$ and mean cost $c_i \in \mathbf{R}_+^d$ for $i$ in $1, \ldots, m$
- **Knapsack:** Known total budget $B \in \mathbf{R}_+^d$ and the length of time horizon $T$
- **Rule:** The decision maker selects one arm to pull at each step and then observes the realized reward and cost of the chosen arm
- **Goal:** Maximize the total reward until any type of budgets is depleted or time up.

## 2. Underlying Linear Programming Problem

- The optimal objective value $\mathrm{OPT}_{\mathrm{LP}}$ of the following LP is an upper bound of the expected optimal reward,

$$\mathrm{OPT}_{\mathrm{LP}} := \max_{x} \ \mu^\top x$$
$$\text{s.t.} \quad Cx \leq B$$
$$x \geq 0$$

where $\mu = (\mu_1, \ldots, \mu_m)^\top$, $C = (c_1, \ldots, c_m)$, and $x = (x_1, \ldots, x_m)^\top \in \mathbf{R}^m$. Denote $x^*$ as its optimal solution and $\mathrm{OPT}_{\mathrm{LP}}$ as its optimal objective value

- If the model is deterministic, the solution of the above LP gives the optimal policy. Here, $x_i$ represents the number of drawing $i$-th arm.
- Denote $\mathcal{I}^* = \{i: x_i^* > 0\}$ as the index set of optimal arms, $\mathcal{J}^* = \{j: (B > Cx^*)_j\}$ as the index set of binding resources, and $\mathcal{I}'$ and $\mathcal{J}'$ as index sets of sub-optimal arms and non-binding resources, respectively.

## 3. Motivation- A Regret Upper Bound

- The regret can be bounded by
$\sum_{i \in \mathcal{J}'} \Delta_i \mathbf{E}[n_i] + \mathbf{E}[B^{(\tau)}]^\top y^*$, where $\Delta_i$ is the reduced cost for the $i$-th arm, $n_i$ is the number of times that $i$-th arm is pulled, $B^{(\tau)}$ is the remaining resources at the termination time, and $y^*$ is the optimal dual price.
- The first term is interpreted as the cost of playing sub-optimal arms; the second term is interpreted as the cost of wasted binding resources.

Thus, to maximize the reward, the decision maker should
- play less sub-optimal arms (optimal arms identification)
- fully consume binding resources (binding resources identification and adaptively procedure to play optimal arms)

## 4. Symmetry between Arms and Knapsack

- With mild conditions, $|\mathcal{I}^*| = |\mathcal{J}^*|$.
- Denote $\mathrm{OPT}_i$ as the optimal objective value of following LP,  $\mathrm{OPT}_i := \max_{x} \ \mu^\top x$,
$$\text{s.t.} \quad Cx \leq B,$$
$$x_i = 0, x \geq 0.$$
then, $i \in \mathcal{I}^* \Leftrightarrow \mathrm{OPT}_i = \mathrm{OPT}_{\mathrm{LP}}$ and $i \in \mathcal{I}' \Leftrightarrow \mathrm{OPT}_i < \mathrm{OPT}_{\mathrm{LP}}$.
- Denote $\mathrm{OPT}_j$ as the optimal objective value of following LP, $\mathrm{OPT}_j = \max_{x} \ \mu^\top x - (B - Cx)_j$,
$$\text{s.t.} \quad \sum_{i=1}^{m} Cx \leq B,$$
$$x \geq 0.$$
Then, $j \in \mathcal{J}^* \Leftrightarrow \mathrm{OPT}_j = \mathrm{OPT}_{\mathrm{LP}}$ and $j \in \mathcal{J}' \Leftrightarrow \mathrm{OPT}_j < \mathrm{OPT}_{\mathrm{LP}}$.
- Set $\delta = \frac{1}{T}(\mathrm{OPT}_{\mathrm{LP}} - \max\{\max_{i \in \mathcal{I}'}\{\mathrm{OPT}_i\}, \max_{j \in \mathcal{J}'}\{\mathrm{OPT}_j\}\})$.
It characterizes the hardness of distinguishing optimal arms from non-optimal arms, and binding resources from non-binding resources.
- $\delta$ is also a generalization of the sub-optimality measure for multi-armed bandit problem.

## 5. UCB and LCB

- When the sample size is not large enough, the bias between real mean and sample mean might be large. This difference might mislead the algorithm to find true optimal arms and binding resources. To avoid this case, we apply the upper confidence bound and lower confidence bound techniques by considering both UCB and LCB of rewards and costs of each arm instead of sample mean as literatures.
- Denote $\mu^U, \mu^L$ as the UCB and LCB of $\mu$, and denote $C^U, C^L$ as the UCB and LCB of $C$. The optimal objective value of following two LPs are UCB and LCB of $\mathrm{OPT}_{\mathrm{LP}}$.

$$\mathrm{OPT}_{\mathrm{LP}}^U := \max_{x} \ (\mu^U)^\top x,$$
$$\text{s.t.} \quad C^L x \leq B,$$
$$x \geq 0.$$

$$\mathrm{OPT}_{\mathrm{LP}}^L := \max_{x} \ (\mu^L)^\top x,$$
$$\text{s.t.} \quad C^U x \leq B,$$
$$x \geq 0.$$

- $\mathrm{OPT}_{\mathrm{LP}}^U$ and $\mathrm{OPT}_{\mathrm{LP}}^L$ converge to $\mathrm{OPT}_{\mathrm{LP}}$
- Similarly, define $\mathrm{OPT}_i^U$, $\mathrm{OPT}_j^U$, $\mathrm{OPT}_i^L$ and $\mathrm{OPT}_j^L$ for all arms and resources. They also satisfy the convergence property.
- As the size of samples are large enough, following inequalities hold: $\mathrm{OPT}_{\mathrm{LP}}^L > \mathrm{OPT}_i^U$ for all optimal-arms and $\mathrm{OPT}_{\mathrm{LP}}^L > \mathrm{OPT}_j^U$ for all non-binding resources. However, $\mathrm{OPT}_{\mathrm{LP}}^L \leq \mathrm{OPT}_i^U$ and $\mathrm{OPT}_{\mathrm{LP}}^L \leq \mathrm{OPT}_j^U$ always hold for all sub-optimal arms and non-binding resources.

## 7. Theoretical Analysis

- During the Phase one, each arm will be played for no more than $O(\frac{\log T}{\delta^2})$ times. Moreover, with probability no less than $1 - O\left(\frac{1}{T^2}\right)$, the algorithm can identify true optimal arms and binding resources.

- If the budget is large enough, $\mathbf{E}[B^{(\tau)}] = O(\frac{1}{\delta^2})$

- The regret of this algorithm can be bounded by $O(\frac{\log T}{\delta^2})$

## 6. Primal-dual Adaptive Algorithm

**Algorithm 1** Primal-dual Adaptive Algorithm for BwK

```
1:  Input: Resource capacity B, T
2:  %% Phase I: Identification of Î* and Ĵ'
3:  Initialize Î* = Ĵ' = ∅, t = 0
4:  Initialize the knapsack process B⁽⁰⁾ = B
5:  while |Î*| + |Ĵ'| < d do
6:      Play each arm i ∈ [m] once
7:      Update t = t + m and the knapsack process B⁽ᵗ⁾
8:      Update the estimates μ̂(t) and Ĉ(t)
9:      Solve the LCB problem and obtain OPT_LP^L(t)
10:     for i ∉ Î* do
11:         Solve the following UCB problem for OPT_i
                OPT_i^U(t) := max_x  (μ^U(t))^⊤ x,
                s.t. C^L(t)x ≤ B,
                     x_i = 0, x ≥ 0.
12:         if OPT_LP^L(t) > OPT_i^U(t) then
13:             Update Î* = Î* ∪ {i}
14:         end if
15:     end for
16:     for j ∉ Ĵ' do
17:         Solve the following UCB problem for OPT_j
                OPT_j^U(t) := min_y  B^⊤y - B,
                s.t. (C^L(t))^⊤ y ≥ μ^U(t) + C_{j,.}^U(t),
                     y ≥ 0.
18:         if OPT_LP^L(t) > OPT_j^U(t) then
19:             Update Ĵ' = Ĵ' ∪ {j}
20:         end if
21:     end for
22: end while
23: Update t = t + 1
24: %% Phase II: Exhausting the binding resources
25: while t < τ do
26:     Solve the following LP
                max_x  (μ^U(t-1))^⊤ x,
                s.t. C^L(t-1)x ≤ B^(t-1),
                     x_i = 0, i ∉ Î*,
                     x ≥ 0.
27:     Denote its optimal solution as x̃
28:     Normalize x̃ into a probability and randomly play an arm according to the probability
29:     Update estimates μ̂(t), Ĉ(t), and B⁽ᵗ⁾
30:     Update t = t + 1
31: end while
```