

Yutong Yan¹, Audrey Durand², Joelle Pineau¹ ¹School of Computer Science, McGill University ²Computer Science and ECSE Department, Université Laval

Abstract

- Multi-armed bandit (MAB) problem is a classic Reinforcement Learning problem where a player faces multiple arms, each associated with a probability distribution over possible rewards [1].
- **Upper Confidence bound applied to trees (UCT)** algorithm has been popular in solving board games. UCT is naturally a better fit to solve Multivariate bandit problems than other approaches by treating the decision-making process as a bandit algorithm for tree search.
- Why to use Bandit Algorithms for factorial **experiments?** Multi-armed bandits minimize the opportunity cost of running an experiment. Previous work has shown Linear Thompson Sampling (LinTS) performs well than traditional heuristics [2].

Factorial Experiment

- Goal: identify the optimal sequence of choices
- Factorial design is composed of factors and choices per factor
 - Each node stores the **history** of previous choices, and each edge represents a **decision**.
 - The order of factors is predefined by the problem • i.e. Graphical design for mobile applications
 - i.e. Adaptive health intervention optimization



Figure 1: factorial design for graphical design with 2 factors x 2 choices per factor

- Typically, there are **many factors** such as:
- Gender, genotype, diet, and experimental protocols
 - Influence the outcome of the experiment
 - Determine the generality of a response.
- Traditional way: separate experiments (A/B testing) in each choice of one factor (very wasteful)
- To include several factors can **avoid** excessive number of experimental subjects.
- Detect interactions amongst intervention components.

Bandit Algorithms for Factorial Experiments



$$\sum_{t=1}^{I} [\mu_{a_*} - \mu_{a_t}]$$

$$B_{a,t} = \hat{\mu}_{a,t} + \sqrt{\frac{2\ln t}{N_{a,t}}}$$

• Linear bandit (LB)

$$x_* = \operatorname{argmax} < \theta_*, x > = \operatorname{argmax} \mu_x$$

$$x \in X_t$$

Action selection policy

$$x_t = \operatorname{argmax} < \hat{\theta}, x > \hat{\theta}$$

$$x \in \mathbf{X}_t$$

$$r_t - < \sigma_*, x_t >$$

ation in Factorial experiments

Bandit Algorithm for Tree Search (BATS)

variance

