

From 6235149080811616882909238708 to 29: Vanilla Thompson Sampling Revisited



Bingshan Hu (University of British Columbia)
Tianyue H. Zhang (Mila-Quebec AI Institute, Université de Montréal)

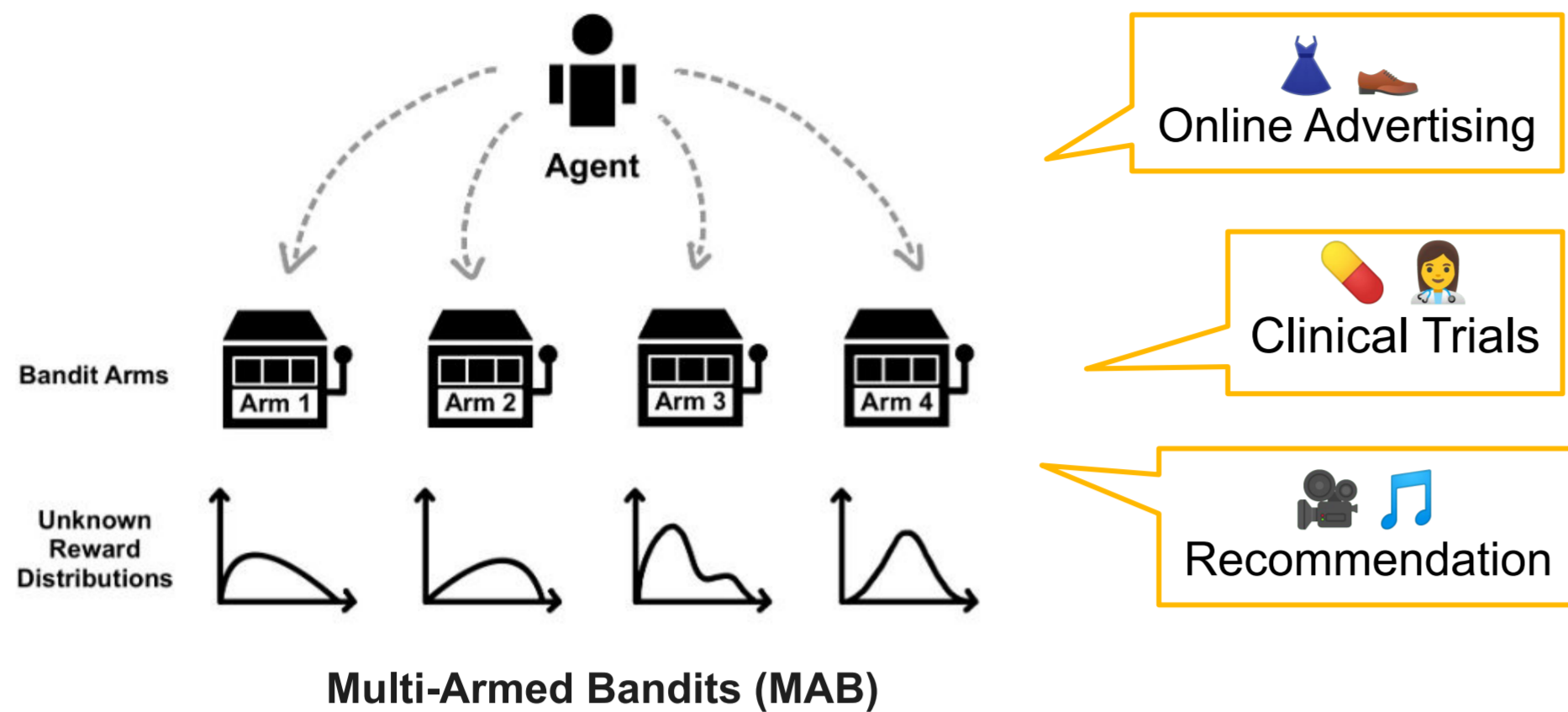


THE UNIVERSITY
OF BRITISH COLUMBIA



Université
de Montréal

Motivation



Stochastic Multi-Armed Bandits

A stochastic MAB instance $\Theta := ([K]; \mu_1, \mu_2, \dots, \mu_K)$
In every round $t = 1, 2, \dots, T$

- Environment generates a reward vector $(X_1(t), \dots, \underbrace{X_j(t)}_{\sim \text{Ber}(\mu_j)}, \dots, X_K(t))$
- Simultaneously, Learner pulls an arm $J_t \in [K]$
- Environment reveals $X_{J_t}(t)$; Learner observes and obtains $X_{J_t}(t)$

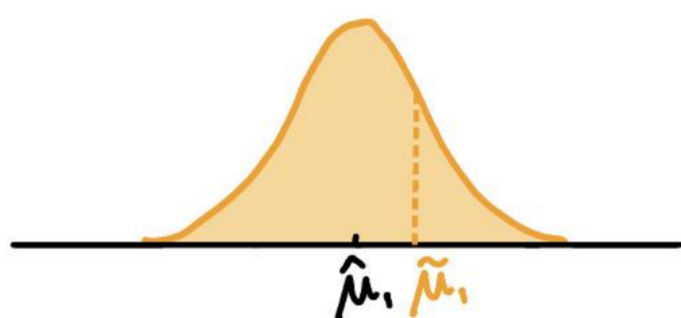
Goal: pull arms sequentially to maximize cumulative reward

$$\text{Regret: } \mathcal{R}(T; \Theta) = \mathbb{E} \left[\sum_{t=1}^T \left(\max_{j \in [K]} \mu_j - \mu_{J_t} \right) \right]$$

Vanilla Thompson Sampling

“Randomly take action according to the probability you believe it is the optimal action” - Thompson 1933

TS uses a **data-dependent distribution** to model the mean of the reward distribution for each arm.



Vanilla TS uses Gaussian distributions to model the mean reward:

- Compute the **empirical mean** of each arm and build the posterior distribution;
- Draw a **random sample** as a proxy for goodness of arm.

Algorithm 1 Thompson Sampling with Gaussian Priors [1]

- Initialization:** for each $i \in [K]$: pull it once to initialize n_i and the empirical mean $\hat{\mu}_{i, n_i}$
- for** $t = K + 1, K + 2, \dots$ **do**
- Draw $\theta_i(t) \sim \mathcal{N}(\hat{\mu}_{i, n_i}, \frac{1}{n_i})$ for all $i \in [K]$
- Pull arm $i_t \in \arg \max_{i \in [K]} \theta_i(t)$ and observe $X_{i_t}(t)$
- Set $n_{i_t} \leftarrow n_{i_t} + 1$ and update the empirical mean $\hat{\mu}_{i_t, n_{i_t}}$ of the pulled arm i_t accordingly.
- end for**

Existing Regret Bound of Vanilla TS

[Agrawal and Goyal, 2017]

$$\sum_{i: \Delta_i > 0} \frac{288 (e^{64} + 6) \ln(T \Delta_i^2 + e^{32})}{\Delta_i} + \frac{10.5}{\Delta_i} + \Delta_i$$

$\Delta_i = \mu_* - \mu_i$
Sub-optimality gap

- The coefficient for the leading term is at least

$$288 \cdot e^{64} \approx 1.8 \times 10^{30}$$

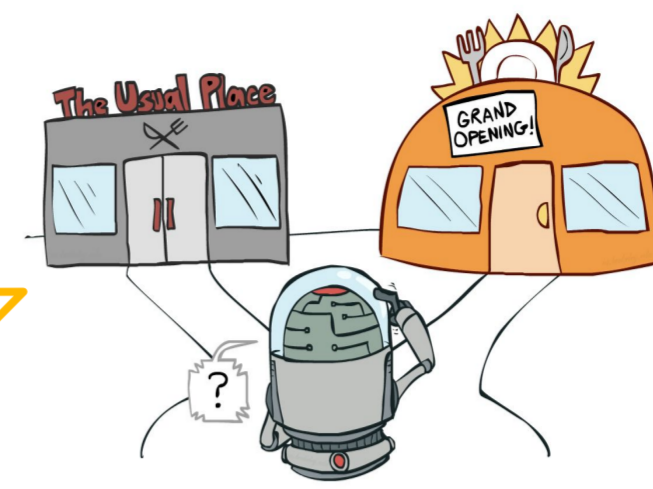
- Since the regret is at most T , the regret bound is vacuous for learning problems when

$$T \leq 288 \cdot e^{64}$$

Challenge: Exploitation vs Exploration Trade-Off

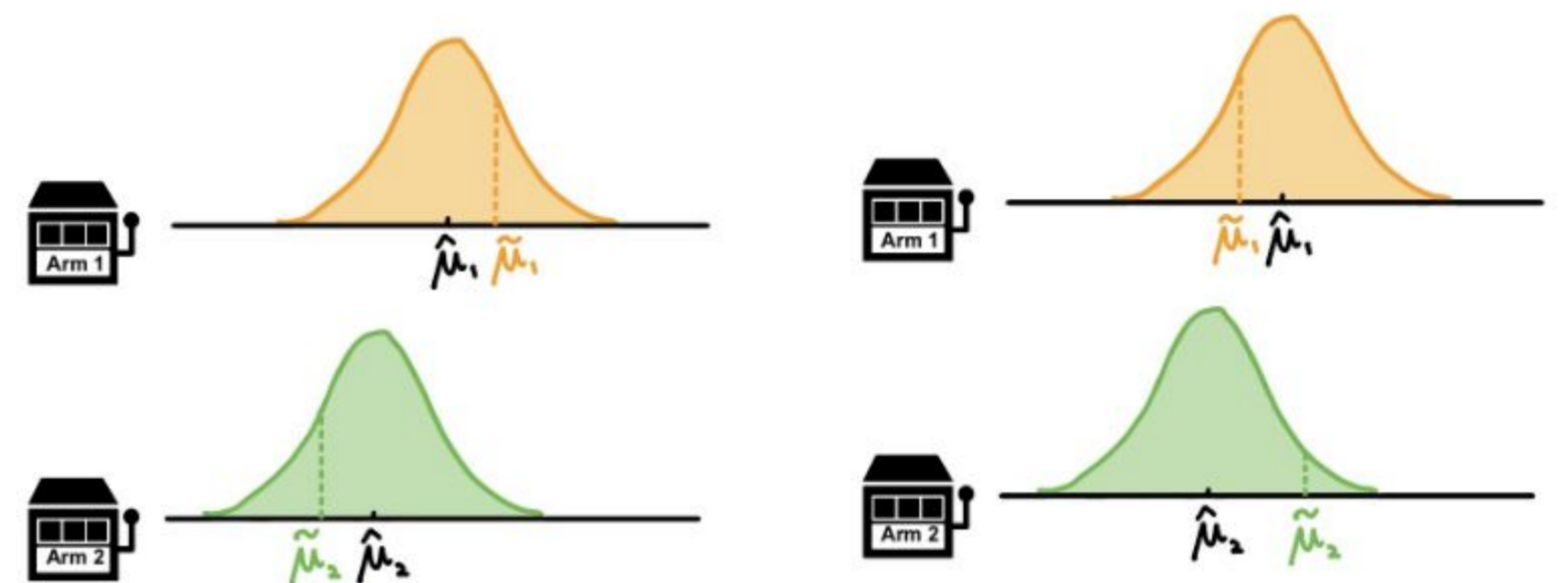
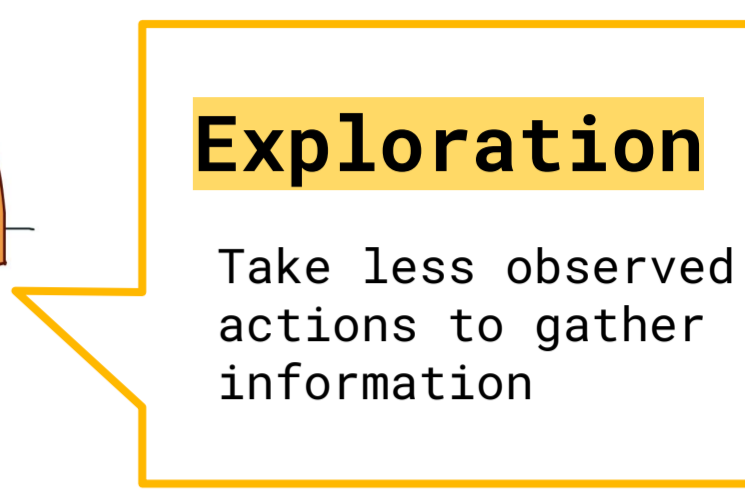
Exploitation

Take actions with high empirical reward to gain pay-off



Exploration

Take less observed actions to gather information

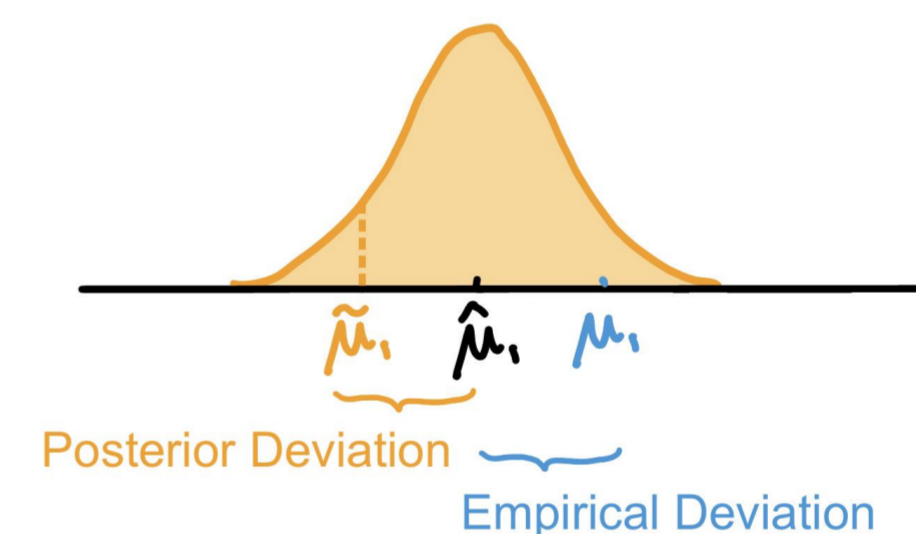


Exploitation

Exploration

Objective

When the posterior distribution of the optimal arm is not concentrated, that is, the optimal arm has not been sufficiently observed,
What is the the expected number of rounds needed before the optimal arm has a good posterior sample?



Example when the **true mean** of the optimal arm is **underestimated**, and the **sample** is also “bad”.

Our Improved Bound

First, we show that the expected number of rounds is at most 29 for us to have a good sample for the optimal arm

Lemma 2 Let $\tau_s^{(1)}$ be the round when the s -th pull of the optimal arm 1 occurs and $\theta_{1,s} \sim \mathcal{N}(\hat{\mu}_{1,s}, \frac{1}{s})$. Then, for any integer $s \geq 1$, we have

$$\mathbb{E}_{\mathcal{F}_{\tau_s^{(1)}}} \left[\frac{1}{\mathbb{P}\{\theta_{1,s} > \mu_1 - \frac{\Delta_1}{2} | \mathcal{F}_{\tau_s^{(1)}}\}} - 1 \right] \leq 29$$

Also, for any integer $s \geq L_{1,i} := \frac{4(\sqrt{2} + \sqrt{3.5})^2 \ln(T \Delta_i^2 + 100^{\frac{1}{3}})}{\Delta_i^2}$, we have

$$\mathbb{E}_{\mathcal{F}_{\tau_s^{(1)}}} \left[\frac{1}{\mathbb{P}\{\theta_{1,s} > \mu_1 - \frac{\Delta_1}{2} | \mathcal{F}_{\tau_s^{(1)}}\}} - 1 \right] \leq \frac{180}{T \Delta_i^2}$$

Then, our improved bound is:

$$\sum_{i: \Delta_i > 0} \frac{1252 \ln(T \Delta_i^2 + 100^{\frac{1}{3}})}{\Delta_i} + \frac{18 \ln(T \Delta_i^2)}{\Delta_i} + \frac{182.5}{\Delta_i} + \Delta_i$$

Note that our improved problem-dependent regret bound also implies an improved worst-case regret bound.

Acknowledgement

This work was supported by Alberta Machine Intelligence Institute (Amii), the Canada CIFAR AI Program and the Natural Sciences and Engineering Research Council of Canada (NSERC) Discovery Grants. Part of the work was done when Bingshan Hu was at Amii and University of Alberta, and when Tianyue Zhang was at University of British Columbia.