



Workshop 10

COMP90051 Statistical Machine Learning
Semester 2, 2024

Learning Outcomes

By the end of this workshop you should be able to:

1. Be able to implement **epsilon-greedy multi-armed bandits**
2. Be able to implement **upper confidence bound multi-armed bandits**
3. Be familiar with offline evaluation of MABs
4. Develop intuition about **exploitation vs. exploration**

Stochastic MAB setting

- Possible actions $\{1, \dots, k\}$ called “**arms**”
 - * Arm i has distribution P_i on bounded **rewards** with mean μ_i
- In round $t = 1 \dots T$
 - * Play action $i_t \in \{1, \dots, k\}$ (possibly randomly)
 - * Receive reward $R_{i_t}(t) \sim P_{i_t}$

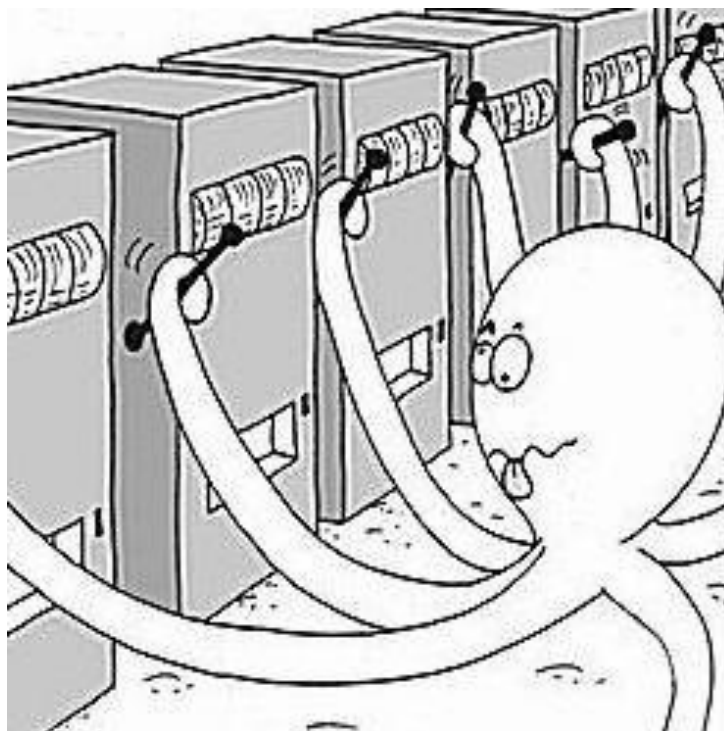
dist $\mu_i = E[R_i]$
 $\rightarrow i_t = \text{arm chosen at time } t$
- Goal: minimise cumulative **regret**
 - * $\mu^* T - \sum_{t=1}^T E[R_{i_t}(t)]$

$\mu^* T - \sum_{t=1}^T \hat{R}_{i_t}(t) \rightarrow$ observed reward
 $\mu^* T$ \rightarrow Expected cumulative reward of bandit
 $\sum_{t=1}^T E[R_{i_t}(t)]$ \rightarrow Best expected cumulative reward with hindsight

where $\mu^* = \max_i \mu_i$
- Intuition: Do as well as a rule that is simple but has knowledge of the future

Multi-armed bandits

each $\xrightarrow{\text{distribution}}$ action has its reward



Open the mab/index.html for a simulation

$$N_{t,k} = \sum_{I=1}^t \mathbb{1}[a_I = k] \quad \text{count } k \text{ was chosen how many times}$$

k : arm number round t .

ϵ -Greedy exploit
 : choose best arm with prob $1-\epsilon$ otherwise

Worksheet 10

μ over n pulls

$$\hat{\mu} = \frac{n\mu + \underbrace{M_t}_{\text{new reward we observe}}}{n+1}$$