# Multi Armed Bandits

DOROZHKO Anton

Novosibirsk State University

May 30, 2020

# Outline

# Framework

MAB is one of the frameworks for **algorithms that make decisions over time under uncertainty**
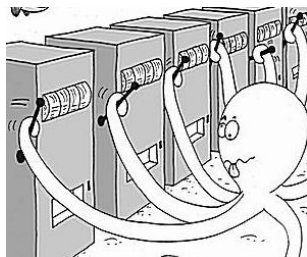
# Examples

1. **News website** : a new user arrives, website picks an article to show, observes user clicks. Goal: maximize the total number of clicks

2. **Dynamic pricing** : an app store, customer arrives, the store chooses the price, the customer buys or leaves forever. Goal: maximize the total profit

3. **Investment** : each morning choose one stock to invest $ . In the end of the day, observe the change in value for each stock. Goal: maximize the total wealth

Novosibirsk
State
University
*THE REAL SCIENCE

# Framework

MAB unifies these examples.
Basic version:

- K possible actions, a.k.a **arms** at each time
- T rounds

# Connection to MDP

### Definition

A Markov Decision Process is a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$

- $\mathcal{S}$ is a space of states
- $\mathcal{A}$ is a space of actions
- $\mathcal{P}$ is a transition probability

$$\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$$

- $\mathcal{R}$ is a reward function :

$$\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$$

- $\gamma$ is a discount factor, $\gamma \in [0, 1]$

# Examples MABs

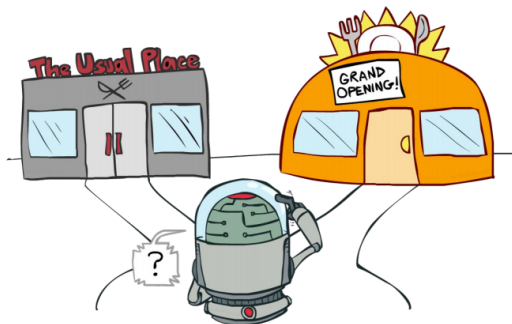| Example | Action | Reward |
|---------|--------|--------|
| News website | an article to display | 1 if clicked, 0 otherwise |
| Dynamic pricing | a price to offer | p is sale, 0 otherwise |
| Ivestment | a stock to invest | change in value during the day |

# Exploration / Exploitation

- observe reward only for chosen arm, not for all
- needs to **explore**
- **explore** = try different arms to get new information
- make optimal neat-term decisions based on available info - **exploitation**

# Exploration vs Exploitation

**TRADEOFF**

learn which arm is the best, but not spend much time learning

# More complex MABs **Feedback**

Auxiliary feedback : other than the reward for chosen arm

| Example | Auxiliary feedback | Reward for any other arm? |
|---|---|---|
| News website | N/A | no |
| Dynamic pricing | sale =>sale at any lower price | yes, for some arms |
| | no sale =>no sale for higher price | |
| Ivestment | change in value for all stocks | yes, for all arms |

- *bandit feedback* : reward for only the chosen arm
- *full feedback* : reward for all arms, that can be chosen
- *partial feedback* : only for some arms

N ✳Novosibirsk
State
University
✳THE REAL SCIENCE

# More complex MABs **Reward**

- *IID rewards* : the reward for each arm is drawn from fixed distribution that depends on the arm, but not on the round t
- *Adversarial rewards*: rewards can be arbitrary, as if they are chosen by "adversary" to fool the agent
- *Constrained adversary*: as Adversarial rewards + some constraints. (e.g. cannot change much from one round to another, ... )
- *Stochastic rewards* : rewards evolves over time as random process, e.g. random walk.

# More complex MAB **Contexts**

Contextual MABs
each round, agent can observe some **context** for each action
goal: learn the best **policy** which maps context to arms, while not
spending much time learning

| Example | Context |
| --- | --- |
| News website | user location and demographics |
| Dynamic pricing | customer's device, location, ... |
| Investments | earning multipliers, state of the company, ... |

# More examples

| Application domain | Action | Reward |
|---|---|---|
| medical trials | which drug to prescribe | health outcom |
| web design | font color or page layout | #clicks |
| content optimization | which item/article to emphasize | #clicks |
| recommender systems | which movie to watch | 1 if follows recommendation |
| datacenter design | which server to route the job to | job completion time |
| robot control | a "strategy" for a given task | job completion time |
| radio networks | which radio frequency to use ? | 1 if successful transmission |
| crowdsourcing | which task to give to which workers, at which price | 1 if task completed at sufficient quality |

# Stochastic Bandits

- Given: K arms, T rounds
- at each round $t \in [T]$
  1. agent picks arm $a_t$
  2. agent observes reawrd $r_t \in [0, 1]$ for the chosen arm
- Goal: maximize total reward over T rounds

# Notation

- arms $a$, rounds $t$
- mean reward of arm $a$ : $\mu(a) = \mathbb{E}[D_a]$
- best mean reward $\mu^* = max_a\mu(a)$
- difference / gap of arm $a$ : $\Delta(a) = \mu^* - \mu(a)$

# Regret : Motivation

- How do we argue if agent is doing a good job ?

# Regret : Motivation

- How do we argue if agent is doing a good job ?
- Different tasks will have different rewards

# Regret : Motivation

- How do we argue if agent is doing a good job ?
- Different tasks will have different rewards
- Some problems have inherently higher rewards

# Regret : Motivation

- How do we argue if agent is doing a good job ?
- Different tasks will have different rewards
- Some problems have inherently higher rewards
- **Standard approach** - compare to the best-arm benchmark
  $\mu^* \cdot T$

# Regret : Definition

### Definition

**Regret at round** $T$ is a difference between the expected reward of always playing and optimal arm and the algorithm's cumulative reward:

$$R(T) = \mu^* \cdot T - \sum_{t=1}^{T} \mu(a_t)$$

# Explore First

1. Exploration phase: try each arm N times
2. Select the arm $\hat{a}$ with the highest average reward (break ties arbitrarily)
3. Exploitation phase: play arm $\hat{a}$ in all the remaining rounds

$N$ is parameter

# Hoeffding's inequality

$$P(|S_n - \mathrm{E}[S_n]| \geq t) \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^{n}(b_i - a_i)^2}\right).$$

# Explore First regret bounds

# Explore First regret bounds 2

# Explore First regret bounds conclusion

Explore-first achieves regret

$$\mathbb{E}[R(T)] \leq T^{2/3} \times O(K \log T)^{1/3}$$

Observations:

- Performance of exploration phase is terrible

# Explore First regret bounds conclusion

Explore-first achieves regret

$$\mathbb{E}[R(T)] \leq T^{2/3} \times O(K \log T)^{1/3}$$

Observations:

- Performance of exploration phase is terrible
- It's better to spread exploration more uniformly over time.

# Explore First regret bounds conclusion

Explore-first achieves regret

$$\mathbb{E}[R(T)] \leq T^{2/3} \times O(K \log T)^{1/3}$$

Observations:

- Performance of exploration phase is terrible
- It's better to spread exploration more uniformly over time.
- E.g. with $\epsilon$-Greedy exploration

# $\epsilon$-Greedy exploration

```
for each round t = 1, 2, . . . do
    Toss a coin with success probability ε_t;
    if success then
        explore: choose an arm uniformly at random
    else
        exploit: choose the arm with the highest average reward so far
end
```

**Algorithm 1.2:** Epsilon-Greedy with exploration probabilities $(\epsilon_1, \epsilon_2, \ldots)$.

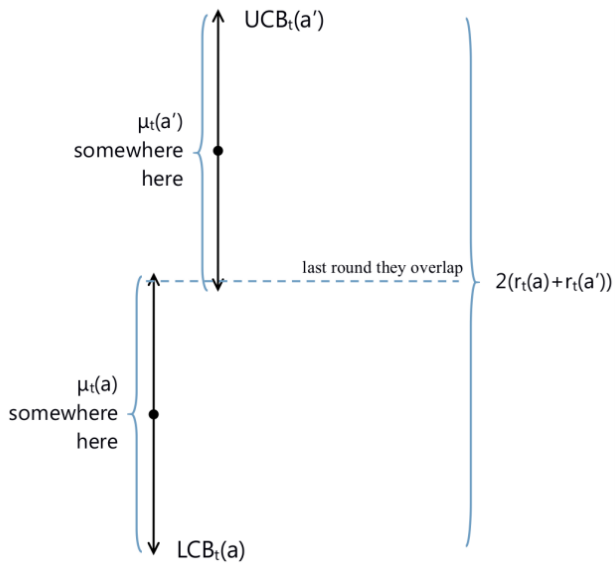# $\epsilon$-Greedy exploration regret

Explore-first achieves regret

$$\mathbb{E}[R(T)] \leq T^{2/3} \times O(K \log T)^{1/3}$$

$\epsilon$-Greedy exploration regret with $\epsilon = t^{-1/3} \cdot (K \log t)^{1/3}$
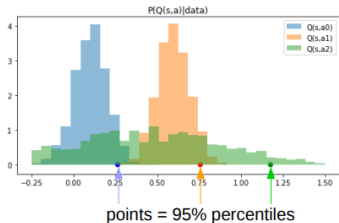
$$\mathbb{E}[R(t)] \leq t^{2/3} \times O(K \log t)^{1/3}$$

for each round t

- $\epsilon$-greedy regret grows linearly
- UCB and Thompson sampling grows with $log(T)$

# Optimism in face of uncertainty

- Policy:
  - Compute 95% upper confidence bound *for each a*
  - Take action with highest confidence bound
  - Adjust: change 95% to more/less



points = 95% percentiles

$Q(s,a)$