

# Multi Armed Bandits

DOROZHKO Anton

Novosibirsk State University

May 30, 2020

# Outline

- 1 Motivation
- 2 Definition
- 3 Regret minimization
- 4  $\epsilon$ -Greedy
- 5 UCB

# Framework

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

## Examples

in MAB

agent? ~~S?~~ A? R?

- 1 News website : a new user arrives, website picks an article to show, observes, user clicks. Goal: maximize the total number of clicks  
*agent* *action*  
*R* *cumulative reward*
- 2 Dynamic pricing : an app store, customer arrives, the store chooses the price, the customer buys or leaves forever. Goal: maximize the total profit  
*A* *agent* *R*  $\{0, 1\}$  *R*  $\{0, p\}$   $\{p_1, p_2\}$
- 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  
*R*  $\{p_1, p_2\}$

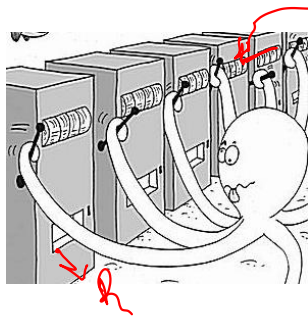
agent? A? R?  
trader choice of stock  $\Delta$

# Framework

MAB unifies these examples.

Basic version:

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



One-armed  
Bandit

# 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~~      *no states / or only 1 MAB*
- $\mathcal{A}$  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	<u>a price to offer</u>	<u>p is sale</u> , 0 otherwise
Investment	a stock to invest	<u>change</u> in value during the day

# Exploration / Exploitation

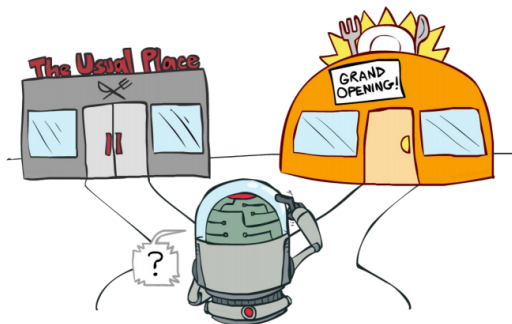
- observe reward only for chosen arm, not for all
- needs to explore
- **explore** = try different arms to get new information
- make optimal near-term decisions based on available info - **exploitation**



# Exploration vs Exploitation

## TRADEOFF

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



<sup>1</sup>CS188 <http://ai.berkeley.edu/home.html>

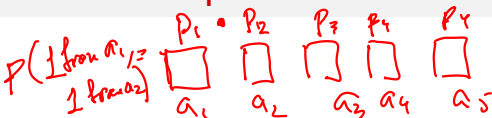
# More complex MABs Feedback

Auxiliary feedback : other than the reward for chosen arm

Example	Auxiliary feedback	Reward for any other arm?
1) News website	N/A	no
2) Dynamic pricing	sale => sale at any lower price no sale => no sale for higher price	yes, for some arms
3) Investment	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

# More complex MABs Reward



identical independent distributions

- 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

$$\mathcal{R}; \pi : S \rightarrow \mathcal{A}$$

## Contextual MABs

$$\text{Contextual Bandits} : \mathcal{X} \rightarrow \mathcal{A}$$

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 outcome <sup>e</sup>
web design	font color or page layout	#clicks <i>A/B</i>
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

Bernoulli bandits

$\xi_t \sim \text{Bernoulli distr}$   
 $p_{a_i} \quad r=1$   
 $1-p_{a_i} \quad r=0$

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

$$\max \sum_{t=1}^T r_t$$

# Notation

*action*

- 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

$$\underbrace{\mu^* \cdot T}$$

$$\mu^* = \max_a \mu(a)$$

# 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:

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

$$\text{expected} \rightarrow \mathbb{E}[R(T)]$$

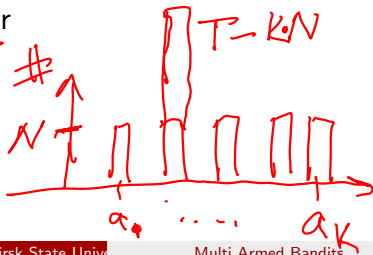
## Explore First

$$1 \quad 2 \quad \dots \quad K$$

for  $K \cdot N$  steps  
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

$$X \sim [a_i, b_i]$$

$$S_n = \frac{X_1 + X_2 + \dots + X_n}{n}$$

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

for Bernoulli  $a_i = 0$   
 $b_i = 1$

$$P(|\bar{X} - E[X]| \geq t) \leq 2e^{-2nt^2}$$

## Explore First regret bounds

$$\mu(a), \mu(a^*)$$

$$P(|x - \mathbb{E}[x]| \geq t) \leq 2e^{-2nt^2}$$

$$\boxed{P_1} \quad \boxed{P_2} \quad \dots \quad \boxed{P_K} \quad \mu(a_i) = p_i \text{ for } \mathcal{B}(p)$$

computes after exploration  
↓

$$P\{| \bar{\mu}(a) - \mu(a) | \leq r(a)\} \geq 1 - 2e^{-2Tr(a)}$$

$$r(a) = \sqrt{\frac{2 \log T}{N}} \quad \geq \underline{1 - \frac{2}{T^4}}$$

suppose "clean event"

$$\underline{\mu(a) - r(a) \leq \bar{\mu}(a) \leq \mu(a) + r(a)}$$

## Explore First regret bounds 2

$$K=2 \quad R(T) = \underbrace{\mu^* \cdot T}_{\text{exploitation}} - \underbrace{\sum \mu(a)}_{\text{exploration}}$$

chosen after  $kN$  steps  $\rightarrow$   $KN$   $\rightarrow$   $T$

best  $a$   $\max \bar{\mu}$

$a \neq a^*$

$$\mu(a) > \mu(a^*) \quad \mu(a) + r(a) \geq \bar{\mu}(a) \geq \bar{\mu}(a^*) - r(a^*)$$

$$\mu(a^*) - \mu(a) \leq r(a) + r(a^*) = O\left(\sqrt{\frac{\log T}{N}}\right) \text{ after } KN$$

$$R(T) = \sum_{t=1}^T (\mu^* - \mu(a_t)) \leq \underbrace{N}_{2N} + \underbrace{O\left(\sqrt{\frac{\log T}{N}}\right)}_{\text{exploit}} \cdot \underbrace{(T-2N)}_{\text{explore}}$$



## Explore First regret bounds conclusion

$$R(T) \leq \uparrow N + O\left(\sqrt{\frac{\log T}{N}} \times T\right) \leftarrow \min R(T) \text{ w.r.t } N$$

Explore-first achieves regret

$$N = \sqrt{\frac{\log T}{N}} T \rightarrow N^3 = \log T \cdot T^2$$

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

$$N = T^{2/3} (\log T)^{1/3}$$

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

played all actions  $N$  times each  $\rightarrow a^*$

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

$1-\epsilon$   $a_t^*$   
 $\epsilon$  random  $a$

# $\epsilon$ -Greedy exploration

```
for each round  $t = 1, 2, \dots$  do  
  Toss a coin with success probability  $\epsilon_t$ ;  
  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  $(\epsilon_1, \epsilon_2, \dots)$ .

# $\epsilon$ -Greedy exploration regret

Explore-first achieves regret

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

all  
# rounds

T

$\epsilon$ -Greedy exploration regret with  $\epsilon = t^{-1/3} \cdot (K \log t)^{1/3} \leftarrow$

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

t

for each round t

- $\epsilon$ -greedy regret grows linearly  $T^{2/3} \cdot 0$
- UCB and Thompson sampling grows with  $\log(T)$

upper confidence bound

Hoeffding's ineq

$$P(|\mu - \mu^*| \leq r) = 1 - \frac{2}{T^r}$$

$$r = \sqrt{\frac{2 \log T}{n_1}}$$

$$r = \sqrt{\frac{2 \log T}{N}}$$

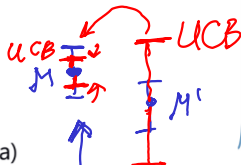
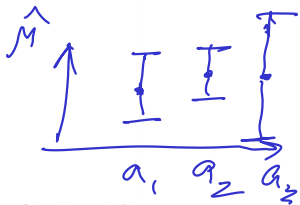
$\mu_t(a)$   
somewhere  
here

$\mu_t(a')$   
somewhere  
here

UCB<sub>t</sub>(a')

last round they overlap

$$\bar{\mu} = \frac{\sum_{i=1}^n r_i}{n}$$

LCB<sub>t</sub>(a) $n_1, n_2, n_3, \dots$ 

$$2(r_t(a) + r_t(a'))$$

$t \nearrow$     $r \nearrow$

$n_1 \nearrow$     $r \searrow$

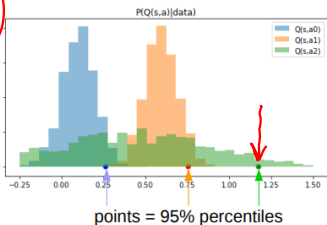
Novosibirsk  
State  
University

\*THE REAL SCIENCE

# Optimism in face of uncertainty

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

*argmax<sub>a</sub> Q(s,a)*  
*ε greedy*



*Q(s, a) → scalar*  
*→ M*

Q(s,a)



Novosibirsk  
State  
University  
\*THE REAL SCIENCE