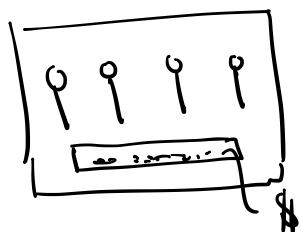


# Multi-armed Bandits

# MAB

- a gambling machine



which one to pull?

- don't know which one is the best
- need to explore arms and exploit good ones

- at  $t$ , when I play arm  $j$ , get reward  $X_j(t)$

↑  
drawn i.i.d. from unknown distn with mean  
reward  $\mu_j$  (assume  $0 \leq X_j(t) \leq 1$ )

- suffer expected regret  $\Delta_j = \mu^* - \mu_j$

$\mu^*$   
mean reward  
of best arm

$\mu_j$   
mean reward for  
arm I chose

- my current estimate of  $\mu_j$ :

$$\hat{\mu}_{j,t} = \frac{1}{T_j(t-1)} \sum_{s=1}^{T_j(t-1)} X_j(s)$$

# times I played j before t

- total regret for whole game

$$R_n = \sum_{t=1}^n \sum_{j=1}^m \Delta_j \mathbf{1}_{\{I_t=j\}}$$

↑ rounds       $\underbrace{\Delta_j}_{\text{regret is } \Delta_j \text{ when we choose arm } j \text{ at time } t}$

## Alg 1: " $\epsilon$ -greedy"

Input:  $n$  rounds

$m$  arms

$k$  constant,  $K > 10$ ,  $K > \frac{4}{\min_j \Delta_j^2} \leftarrow$  difference between  
best arm and 2nd best arm

$$\{\epsilon_t\}_{t=1}^n = \min \left\{ 1, \frac{km}{t} \right\} \longleftarrow \frac{km}{t}$$

Initialize: play all arms once.  $\hat{X}_{j,t} = X_j(j)^t \quad j=1..m$

while  $t \leq n$

choose arm  $j$   
with prob  $1/m$

with prob  $\varepsilon_t$  play an arm uniformly at random "explore"  $\leftarrow$

otherwise (w. prob  $1-\varepsilon_t$ ) play best arm, play  $j$  s.t.  $\hat{X}_{jt} \geq \hat{X}_{it} \forall i$  "exploit"

get reward  $X_j(t)$

update  $\hat{X}_{jt}$

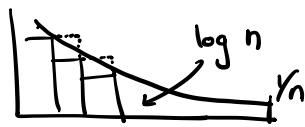
end

Regret bounds are the theory behind MAB algs.

"Theorem 1"  $\mathbb{E}[R_n] \leq O(\log n)$

$$\text{Note : } \sum_{t=1}^n \frac{1}{t} \leq \log n + 1$$

$$\sum_{t=1}^n \frac{1}{t} > \int_1^{n+1} \frac{1}{t} dt = \ln(n+1)$$



Theorem 1 (Auer, Cesa-Bianchi, Fischer)

$$\begin{aligned}
 F(R_n) &\leq \sum_{j=1}^m \Delta_j \quad \leftarrow \text{initialization} \\
 &+ \sum_{t=m+1}^n \sum_{j: \mu_j < \mu_*} \Delta_j \left( \varepsilon_t \frac{1}{m} + (1-\varepsilon_t) \beta_j(t) \right) \\
 &\quad \text{where } \beta_j(t) = k \left( \frac{t}{m k e} \right)^{-\frac{1}{k-1}} \log \left( \frac{t}{m k e} \right) + \frac{4}{\Delta_j} \left( \frac{t}{m k e} \right)^{-\frac{1}{k-1}} - \frac{k \Delta_j^2}{4} \\
 &\quad \text{decays faster than } \frac{1}{t} \quad o(\frac{1}{t})
 \end{aligned}$$

Proof idea :

$$\Delta_j (\varepsilon_t) = \underbrace{m}_{\substack{\text{prob to play} \\ \text{arm } j}} + \underbrace{(1-\varepsilon_t)}_{\substack{\text{prob to explore}}} P(\hat{X}_{j,T_j(t-1)}) \geq \hat{X}_{i,T_i(t-1)} + v_i)$$

↑  
prob to play  
arm  $j$  when exploiting.  
Pr you think arm  $j$  is  
the best when it's not

↑  
part of proof is to  
upper bound this by  $B$

Proof: if we think arm  $j$  is the best  $\alpha \leftarrow c$  these must have happened:

- 1) we overestimated  $M_*$  by a lot or
  - 2) we underestimated  $M_*$  by a lot

$$\begin{array}{c}
 \text{or} \\
 \Delta_{j/2} \left\{ \begin{array}{l} -\mu_x \\ -\mu_j \end{array} \right. \quad \begin{array}{l} -\hat{x}_j \\ -\hat{x}_7 \\ -\hat{x}_* \\ -\hat{x}_2 \end{array} \\
 \Delta_{j/2} \left\{ \begin{array}{l} -\mu_x \\ -\mu_j \end{array} \right. \quad \vdots
 \end{array}$$

$$\hat{X}_{j, T_j(t-1)} \geq \mu_j + \frac{\Delta_j}{2}$$

$$\left. \begin{array}{c} \mu_x \\ \mu_s/2 \\ \mu_j/2 \\ \mu_z \end{array} \right\}$$

$$\text{Step 1: } P(\hat{X}_{j, T_j(t-1)} \geq \hat{X}_{i, T_i(t-1)} \forall i) \leq \underbrace{P(\hat{X}_{j, T_j(t-1)} \geq \mu_j + \frac{\Delta_j}{2})}$$

$$P(A \cup B) = P(A) + P(B)$$

## Step 2 : Clever Bounding

$$\textcircled{*1} = P(\hat{X}_{j,T_j(t-1)} \geq \mu_j + \frac{\Delta_j}{2}) = \sum_{s=1}^{t-1} P(T_j(t-1) = s, \hat{X}_{j,s} \geq \mu_j + \frac{\Delta_j}{2})$$

$$= \underbrace{\sum_{s=1}^{t-1} P(T_j(t-1) = s | \hat{X}_{j,s} \geq \mu_j + \frac{\Delta_j}{2})}_{e^{-\frac{\Delta_j^2}{2}s}} \underbrace{P(\hat{X}_{j,s} \geq \mu_j + \frac{\Delta_j}{2})}_{\text{the fading}}$$

$$= \underbrace{\sum_{s=1}^{t-1} P(T_j(t-1) = s \mid \hat{X}_{j,s} \geq \mu_j + \frac{\Delta_j}{2})}_{\text{the fading}} \underbrace{P(\hat{X}_{j,t} \geq \mu_j + \frac{\Delta_j}{2})}_{e^{-\frac{\Delta_j^2}{2}s}}$$

if  $s$  is big  
this term is small

$$x_0 := \frac{1}{2m} \sum_{s'=1}^t \epsilon_{s'} \stackrel{\text{def}}{=} \sum_{s=1}^{\lfloor x_0 \rfloor} " + \sum_{s=\lfloor x_0 \rfloor + 1}^{t-1} "$$

" " " "

$$\stackrel{\text{def}}{=} \left( \frac{1}{\Delta_j^2} \right) e^{-\frac{\Delta_j^2}{2} \lfloor x_0 \rfloor}$$

So far:

$$*\textcircled{1} \leq e^{-\frac{1}{5} x_0} + \frac{2}{\Delta_j^2} e^{-\frac{(\Delta_j/2)}{2} \lfloor x_0 \rfloor}$$

$$x_0 \geq \frac{k}{2} \log \frac{t}{mke}$$

$$*\textcircled{1} \leq \frac{k}{2} \log \left( \frac{t}{mke} \right) \cdot \left( \frac{t}{mke} \right)^{-k/10} + \frac{2}{\Delta_j^2} \left( \frac{t}{mke} \right)^{-k\Delta_j^2/4}$$

Step 3: Turns out, the same bound holds for  $*\textcircled{2}$

Step 4: Combine Steps 1, 2, 3.

$$P(\hat{X}_{j,\tau_j(t-1)} \geq \hat{X}_{i,\tau_i(t-1)} \mid h_i) \leq *\textcircled{2}$$

UCB Alg: create an upper confidence bound on each arm  
pick the arm with the highest UCB

input: n rounds

m arms

initialize: play all arms once, initialize  $\hat{X}_j$  for  $j=1 \dots m$

for  $t = m+1 \dots n$

play arm  $j$  with highest UCB

$$\hat{X}_{j,T_j(t-1)} + \sqrt{\frac{2 \log(t)}{T_j(t-1)}}$$

get reward  $X_j$

update  $\hat{X}_j$

end

I I I - I

Theorem: regret grows logarithmically in n

$$\mathbb{E}[R_n] \leq \sum_{j=1}^m \Delta_j$$

$$+ \sum_{j: \mu_j < \mu_*} \frac{8}{\Delta_j} \log n + \sum_{j=1}^m \Delta_j \left( \sum_{t=m+1}^n 2 t^{-4} (t-1-m)^2 \right)$$