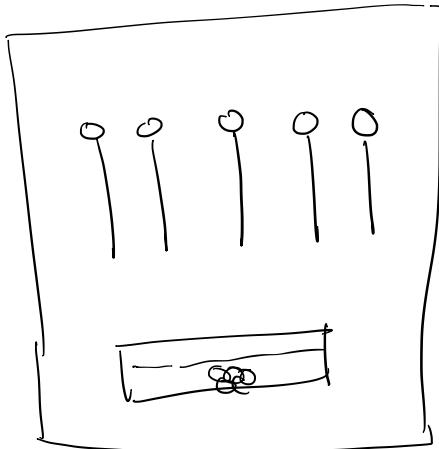


Multi-armed Bandits

MAB

- a gambling machine



- need to explore arms and exploit good ones

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

\uparrow
 drawn iid from unknown dist
 with mean reward μ_j
 (assume $|X_j(s)| \leq 1$)

- suffer expected regret $\Delta_j = \mu_x - \mu_j$

\uparrow mean reward of best arm
 \uparrow mean reward for arm I chose

- my current estimate of μ_j :

$$\hat{\mu}_j = \underbrace{\frac{1}{T_j(t-1)}}_{\# \text{times I played } j \text{ before } t} \sum_{s=1}^{T_j(t-1)} X_j(s)$$

- Total regret for whole game

$$R_n = \sum_{t=1}^n \sum_{j=1}^m \Delta_j \underbrace{\mathbf{1}_{\{I_t=j\}}}_{\substack{\text{regret is } \Delta_j \text{ when} \\ \text{we chose arm } j \text{ at} \\ \text{time } t.}}$$

Alg 1 - "ε-greedy"

Input: n rounds

m arms

$$k \text{ constant}, k > 10 \text{ and } R > \frac{4}{\min_j \Delta_j}$$

$$\{\varepsilon_t\}_{t=1}^n = \min\left\{1, \frac{Rm}{t}\right\} \leftarrow \text{prob to explore}$$



Initialize: play all arms once. $\hat{X}_j = X_j, j=1-m$

while $t < n$

w.p. ε_t play an arm uniformly at random "explore" choose arm j with prob $\frac{1}{m}$

otherwise (w.p. $1-\varepsilon_t$) play best arm, play j s.t. $\hat{X}_j \geq \hat{X}_i \forall i$ "exploit"

get reward X_j

update \hat{X}_j

end

Regret bounds are the theory behind MAB alg.

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

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

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

Theorem 1 (Auer, Cesa-Bianchi, Fischer)

$$\mathbb{E}(R_n) \leq \sum_{j=1}^m \Delta_j \quad \xrightarrow{\text{initialization}}$$

$$+ \sum_{t=m+1}^n \sum_{j: \mu_j < \mu_*} \Delta_j \left(\varepsilon_t \frac{1}{m} + (1-\varepsilon_t) \beta_j(t) \right) \quad \xrightarrow{\text{negative}}$$

$$\text{where } \beta_j(t) = k \left(\frac{t}{mke} \right)^{\frac{1}{k-1}} \log \left(\frac{t}{mke} \right) + \frac{1}{\Delta_j^2} \left(\frac{t}{mke} \right)^{\frac{k-1}{k}}$$

Proof idea:

$$\Delta_j \left(\varepsilon_t \frac{1}{m} + (1-\varepsilon_t) P(\hat{X}_{j, T_j(t-1)} = \hat{X}_{*, T_*(t-1)}, \forall i) \right) \quad \xrightarrow{\text{upper bound this}}$$

Proof: if we think arm j is the best one of these must have happened.

- 1) We overestimated μ_j by a lot or $\hat{X}_j - \mu_j$
 - 2) We underestimated μ_* by a lot or $\hat{X}_* - \mu_*$
- $\hat{X}_{j, T_j(t-1)} \geq \mu_j + \frac{\Delta_j}{2}$
- $\hat{X}_{*, T_*(t-1)} \leq \mu_* - \frac{\Delta_j}{2}$

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

Step 2 : Clever bounding

$$\begin{aligned} \textcircled{1} &= P(\hat{X}_{j,T_j(t-1)} \geq \mu_j + \frac{\Delta_j}{2}) = \sum_{s=1}^{t-1} P(T_j(t-1) = s, \hat{X}_{j,s} \geq \mu_j + \frac{\Delta_j}{2}) \\ &= \sum_{s=1}^{t-1} P(T_j(t-1) = s \mid \hat{X}_{j,s} \geq \mu_j + \frac{\Delta_j}{2}) P(\hat{X}_{j,s} \geq \mu_j + \frac{\Delta_j}{2}) \\ &= \sum_{s=1}^{t-1} P(T_j(t-1) = s \mid \hat{X}_{j,s} \geq \mu_j + \frac{\Delta_j}{2}) \cdot e^{-\frac{\Delta_j^2}{2}s} \xrightarrow{\text{oeffding}} \\ &\quad \uparrow \text{If } s \text{ is small, } \uparrow \text{ could be big...} \\ &\quad \uparrow \text{If } s \text{ is big we're ok thanks to} \end{aligned}$$

$$s \text{ small} \rightarrow S \leq X_o := \frac{1}{2m} \sum_{s=1}^t \varepsilon_s$$

$$\begin{aligned} \textcircled{2} &\downarrow \lfloor X_o \rfloor \\ &= \sum_{s=1}^{\lfloor X_o \rfloor} + \sum_{s=\lfloor X_o \rfloor}^{t-1} \\ &= \sum_{s=1}^{\lfloor X_o \rfloor} P(T_j(t-1) = s \mid \hat{X}_{j,s} \geq \mu_j + \frac{\Delta_j}{2}) \cdot 1 + \sum_{s=\lfloor X_o \rfloor}^{t-1} 1 \cdot e^{-\frac{(\Delta_j^2/2)s}{\Delta_j^2/2}} \\ &\quad \uparrow \text{random } T_j^{(\text{random})}(t-1) = \# \text{ times } j \text{ is chosen during exploration} \\ &\quad \uparrow \lambda \textcircled{3} \end{aligned}$$

$$\sum_{s=1}^{\lfloor X_o \rfloor} P(T_j^{(\text{random})}(t-1) \leq s \mid \hat{X}_{j,s} \geq \mu_j + \frac{\Delta_j}{2})$$

$\Downarrow - T_j^{(\text{random})} \text{ indep of } \hat{X}_{j,s}$

$$\sum_{s=1}^{\lfloor X_o \rfloor} P(T_j^{(\text{random})}(t-1) \leq s)$$

$$\sum_{s=1}^{\lfloor X_o \rfloor} \frac{\lambda}{\lfloor X_o \rfloor} P(T_j^{(\text{random})}(t-1) \leq \lfloor X_o \rfloor)$$

$$\lfloor X_o \rfloor P(T_j^{(\text{random})}(t-1) \leq \lfloor X_o \rfloor)$$

$\Downarrow \text{ skipping steps: apply tail bound}$

$e^{-\gamma_5 \lfloor X_o \rfloor}$

So far,

$$(*1) \leq \lfloor X_0 \rfloor e^{-\gamma_S \lfloor X_0 \rfloor} + \frac{2}{\Delta_j^2} e^{-\left(\frac{\Delta_j^2}{2}\right) \lfloor X_0 \rfloor}$$

A

as $\lfloor X_0 \rfloor$ increases \uparrow decreases. The worst is if
 $\lfloor X_0 \rfloor$ is small.

$$\begin{aligned} X_0 &= \frac{1}{2m} \sum_{s=1}^t \varepsilon_s = \frac{1}{2m} \sum_{s=1}^t \min \left\{ 1, \frac{k_m}{s} \right\} \\ &= \frac{1}{2m} \left[\sum_{s=1}^{\lfloor k_m \rfloor} 1 + \sum_{s=\lfloor k_m \rfloor+1}^t \frac{k_m}{s} \right] \\ &= \frac{\lfloor k_m \rfloor}{2m} + \underbrace{\frac{k_m}{2m} \sum_{s=\lfloor k_m \rfloor+1}^t \frac{1}{s}}_{\stackrel{VI}{=}} \\ &\quad \log(t+1) - \log(\lfloor k_m \rfloor + \log e) \\ &= \frac{k}{2} \frac{\lfloor k_m \rfloor}{mke} + \frac{k}{2} \log \left(\frac{t+1}{\lfloor k_m \rfloor e} \right) \\ &\geq \frac{k}{2} \log \frac{\lfloor k_m \rfloor}{mke} + \frac{k}{2} \log \left(\frac{t}{\lfloor k_m \rfloor e} \right) \\ &= \frac{k}{2} \log \left(\frac{\lfloor k_m \rfloor}{mke} \frac{t}{\lfloor k_m \rfloor e} \right) = \frac{k}{2} \log \frac{t}{mke} \end{aligned}$$

plug into A

$$\begin{aligned} (*1) &\leq \frac{k}{2} \log \left(\frac{t}{mke} \right) e^{-\gamma_S \frac{k}{2} \log \frac{t}{mke}} + \frac{2}{\Delta_j^2} e^{-\frac{\Delta_j^2}{2} \left(\frac{k}{2} \log \frac{t}{mke} \right)} \\ &= \frac{k}{2} \log \left(\frac{t}{mke} \right) \cdot \left(\frac{t}{mke} \right)^{-\frac{k}{10}} + \frac{2}{\Delta_j^2} \left(\frac{t}{mke} \right)^{-\frac{k \Delta_j^2}{4}} \end{aligned}$$

Step 3: Turns out the same bound holds for (*2)

Step 4: Combine Steps 1, 2, 3:

$$P\left(\hat{X}_{j, T_j(t-1)} \geq \hat{X}_{i, T_i(t-1)} + i\right) \leq k \left(\frac{t}{mke} \right)^{-\frac{k}{10}} \log \left(\frac{t}{mke} \right) + \frac{4}{\Delta_j^2} \left(\frac{t}{mke} \right)^{-\frac{k \Delta_j^2}{4}}$$

□

UCB Alg: Create an upper conf bound on each arm.
Pick the arm with the highest UCB

input: n rounds
m arms

initialize: play all arms once, initialize \hat{X}_i

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

Theorem:

$$\mathbb{E}(R_n) \leq \sum_{j=1}^m \Delta_j$$

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

