



Multi-Armed
Bandit

MULTI Armed Bandit Problem

EXPLORE VS EXPLOIT

→ try different machine

→ Pull arm on "best" machine

$\mathcal{Y} = \text{Set of Strategies}$
 $= [n]$

$\mathcal{I} = \text{Set of all possible cost functions}$
Costs in range $[0, 1]$

$\mathcal{I} = [0, 1]^{[n]}$

Sequence of cost functions

$$C_1, C_2, \dots, C_T \in \Gamma$$

we pick

$$x_t \text{ for } t \geq 1$$

depends only on

$$\left\{ C_i(x_i) \right\}_{i=1}^{t-1}$$

$$\hat{R}(\text{ALG}, T) = \max_{x \in \mathcal{X}} \frac{1}{T} \mathbb{E} \left[\sum_{t=1}^T \ell_t(x_t) - \ell_t(x) \right]$$

Adversary

can pick cost functions
ahead of time

OR

can have adaptive Adversary,
can know $x_1 \dots x_{t-1}$ but $x_1 \dots x_T$

Ethical Clinical Trials

typically

$n = \#$ of treatments

$T = \#$ of patients in trial

x_t is the treatment for patient t

$C_t(x_t)$ is how well treatment works

0 = success
1 = Death

Server Selection in networks

we have n servers

x_t is server we pick at DNS query t

$c_t(x_t)$ is the time it took

Internet Ad Placement

first we get a visitor

n ads we could show

x_t is the ad that we show them

$$c_t(x_t) = \begin{cases} 1 & \text{if click ad} \\ 0 & \text{if no click} \end{cases}$$

I identical goods for sale

y = prices you could offer

x_t is the price you offer
visitor t

$c_t(x_t) = 1$ if they buy
or 0 otherwise

Overlay routing in networks

$$S = v_0, v_1, \dots, v_L = t$$

$$d(v, w)$$

$$\text{total time} = \sum_{i=0}^{L-1} d(v_i, v_{i+1})$$

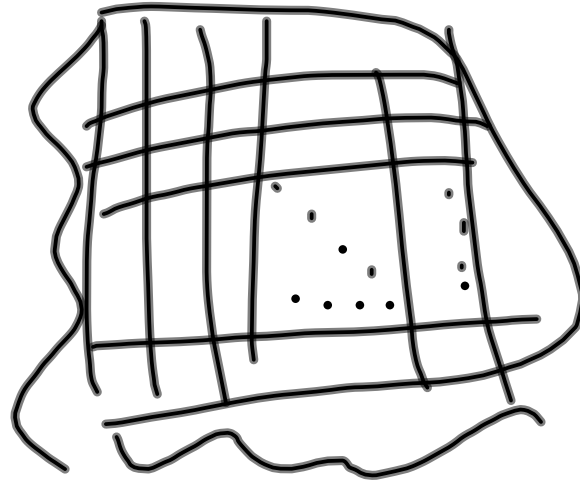
\mathcal{P} = set of paths from S to t

~~$C(x, z) =$~~

GO

S = moves
you could
make

19

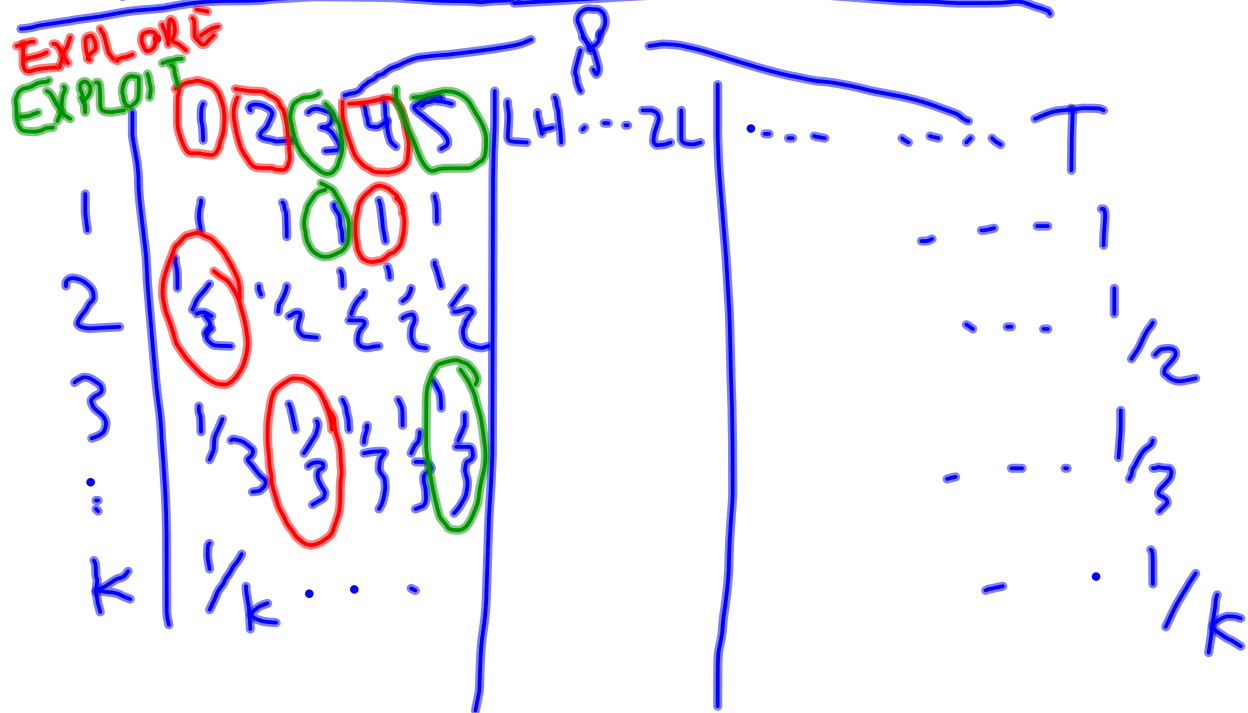


X_t = more than you pick

19

$C_t(x_t)$ = win rate
for that move at that time

Phased Simulation



Pick $L > k (=n)$

(Assume $T = P \cdot L$)

REX depends on parameter $\varepsilon > 0$
Hedge FPL \Rightarrow BEX(ε)

$$\hat{R}(\text{BEX}(\varepsilon), t) \leq \varepsilon + O\left(\frac{\log n}{\varepsilon t}\right)$$

Function $\tau_j(i) \rightarrow \{jL, \dots, (j+1)L\}$

$\tau_j(i)$ is 1-to-1. chosen randomly.

τ_j represents the time steps
in interval j for which we are EXPLORING.

The exploit steps are then the remaining
 $[T] \setminus \text{EXPLORE steps.}$

\hat{C}_j is func. such that $\hat{C}_j(i) = C_{t_j(i)}(i)$

P_j : dist. over $[n]$ from BEX
at time j .

Pick $X_t = \begin{cases} i & \text{if } \tau_j(i) = t \\ \hat{x}_j & \text{otherwise,} \end{cases}$

where $\hat{x}_j \leftarrow p_j$.

Phased Simulation

	EXPLORE					EXPLOIT				...	T	
	1	2	3	4	5	6	7	8	9	10	...	T
1	1	1	1	1	1	1/2	1/2	1/2	1/2	1/2	...	1
2	1/2	1/2	1/2	1/2	1/2	1	1	1	1	1	...	1/2
3	1/3	1/3	1/3	1/3	1/3	1/3	1/3	1/3	1/3	1/3	...	1/3
⋮												
k	1/k	1/k	1/k	1/k	1/k	1/k	1/k	1/k	1/k	1/k	...	1/k

$$\hat{C}_1 = \begin{cases} 1 \rightarrow 1 \\ 2 \rightarrow 1/2 \\ 3 \rightarrow 1/3 \end{cases}$$

Pick $L > k (=n)$
 (Assume $T = P \cdot L$)

$$P_j \leftarrow \begin{cases} 1 \rightarrow 3/4 \\ 2 \rightarrow 3/10 \\ 3 \rightarrow 1/6 \end{cases}$$

Analysis of PSIM = "phased simulation"

$$\frac{1}{P} \sum_j \hat{C}_j \approx \frac{1}{T} \sum_t C_t$$

	1	2	...	L	...	(k-1)L+1	...	kL	...	T	
1	0	0	...	0	...	0	...	0	7/8	...	7/8
2	1	1	...	1	...	1	...	1	1	...	1
3	1	1	...	1	...	1	...	1	1	...	1
...	1	1	...	1	...	1	...	1	1	...	1
$\frac{1}{T} \sum C(2)$	0	0	...	0	...	0	...	0	7/8	...	7/8
$\frac{1}{T} \sum C(1)$	1	1	...	1	...	1	...	1	1	...	1

$\frac{1}{T} \sum C(2) = 1/2$ $\frac{1}{T} \sum C(1) = 7/8$

Assume C_t is not
adaptive.

Lemma 1. $E[\hat{c}_j | \mathcal{F}_{<j}] = \bar{c}_j$.

$$= \frac{1}{L} \sum_{t \in \Phi_j} c_{t,j}$$

Proof: LHS & RHS: $[n] \rightarrow [p, 1]$.

Let arb. $i \in [n]$. WTS:

$$E[\hat{c}_j(i) | \mathcal{F}_{<j}] = \bar{c}_j(i) \checkmark$$

↳ obvious. By unif. dist. of τ_j .
kind of.

Lemma 2. Let $x \in \mathcal{Y}$,

$$E[\hat{c}_j(\hat{x}_j) | \mathcal{F}_{< j}] = E[\bar{c}_j(\bar{x}_j) | \mathcal{F}_{< j}]$$

Proof. $E[\hat{c}_j(\hat{x}_j) | \mathcal{F}_{< j}]$

$$= E\left[\sum_x p_j(x) \hat{c}_j(x) | \mathcal{F}_{< j}\right]$$

$$= \sum_x p_j(x) E[\hat{c}_j(x) | \mathcal{F}_{< j}]$$

$$= E[\bar{c}_j(\bar{x}_j) | \mathcal{F}_{< j}].$$

Lemma 3. Let $x \in \mathcal{Y}$.

$$\frac{1}{T} \mathbb{E} \left[\sum_{t \in \text{EXPLOIT}} (c_t(x_t) - c_t(x)) \right] \leq \zeta + \frac{\log n}{\varepsilon P}.$$

Proof.

$$\frac{1}{T} \sum_{t \in \text{EXPLOIT}} \mathbb{E} [c_t(x_t) - c_t(x)]$$

$$\leq \frac{1}{T} \sum_{t=1}^T \mathbb{E} [c_t(\hat{x}_t) - c_t(x)]$$

$$= \frac{1}{P} \sum_{j=1}^P \mathbb{E} [\bar{c}_j(\hat{x}_j) - \bar{c}_j(x)]$$

$$\leq \sum + \frac{\log n}{\epsilon p}$$

Lemma 4. Let $x \in \mathcal{Y}$.

$$\text{Then } \frac{1}{T} \sum_{t \in \text{EXPLORE}} [c_t(x_t) - c_t(x)] \leq \frac{nP}{T}.$$

Proof. $c_t(x_i) - c_t(x) \leq 1.$

EXPLORE has $n \cdot P$

Theorem 5. Let $\varepsilon = \left(\frac{n \log n}{T}\right)^{1/3}$,

$$p = \log^{1/3}(n) \left(\frac{\varepsilon}{n}\right)^{2/3}.$$

Then, $\hat{R}(PSim, T) = O\left(\left(\frac{n \log n}{T}\right)^{1/3}\right)$

Proof. $\hat{R}(PSim, T) \leq \varepsilon + \frac{\log n}{\varepsilon p} + \frac{n p}{T}.$