

# Multi-Armed Bandit

Given unknown probability distributions  $R_1, \dots, R_k$  with means  $\mu_1, \dots, \mu_k$  ( $\mu^* = \max_{1 \leq i \leq k} \mu_i$ )

Choose indices  $i_1, i_2, \dots$  to optimize payout

Regret = difference between cumulative payout and best expectation

$$R_T = T \cdot \mu^* - \sum_{t=1}^T \hat{r}_t \quad \hat{r}_t = \text{reward at time } t$$

"optimal" means  $P\left(\lim_{T \rightarrow \infty} \frac{R_T}{T} = 0\right) = 1$  zero regret

	Arm 1	Arm 2	Arm 3
Ex:	prob $\frac{1}{3}$	prob $\frac{1}{4}$	prob $\frac{1}{100}$
	payout 2	payout 3	payout 200
	$\frac{2}{3}$	$\frac{1}{4}$	$\frac{59}{100}$
	0	$\frac{1}{2}$	0
	$\mu_1 = \frac{2}{3}$	$\frac{1}{2}$	$\mu_3 = 2 = \mu^*$
		$\mu_2 = \frac{3}{8}$	

uniform rotation: 1 2 3 1 2 3 1 2 3

$\downarrow$  avg regret  $2 - \frac{2}{3} = \frac{4}{3}$   
 $\swarrow$  avg regret  $2 - \frac{3}{8} = \frac{9}{8}$   
 $\leftarrow$  avg regret = 0  
 tot expected regret over seq of 3 plays =  $\frac{4}{3} + \frac{9}{8} = \frac{59}{24}$

$$\lim_{T \rightarrow \infty} \frac{R_T}{T} = \frac{59}{72}$$

greedy: play each once, play arm w/ best result forever

1 2 3 1 1 1 1 1 ...  $\lim_{T \rightarrow \infty} \frac{R_T}{T} = \frac{4}{3}$

1 2 3 2 2 2 2 2 ...  $\lim_{T \rightarrow \infty} \frac{R_T}{T} = \frac{9}{8}$

1 2 3 3 3 3 3 3 ...  $\lim_{T \rightarrow \infty} \frac{R_T}{T} = 0$

} occur with prob > 0

$$P\left(\lim_{T \rightarrow \infty} \frac{R_T}{T} = 0\right) < 1$$

$\epsilon$ -greedy: play each arm once, then play randomly chosen arm w/ probability  $\epsilon$  arm with best avg observed reward otherwise

$$\lim_{T \rightarrow \infty} \frac{R_T}{T} = \frac{\epsilon}{k} \cdot \sum_{i=1}^k (\mu^* - \mu_i) > 0 \quad (\text{assuming not all arms are optimal})$$

zero regret

Choose arm  $j$  that maximizes

UCB  
("upper confidence bound")

$$\bar{r}_j + \sqrt{\frac{2 \ln T}{n_j}}$$

tunable parameter

total plays

avg observed reward for arm  $j$

number of times arm  $j$  played

exploitation

exploration