

## Week 7 - Lecture 32

### 1) Greedy Algorithm

$$I_t = \underset{i}{\operatorname{argmax}} \hat{\mu}_i(t-1)$$

Arm 1 : Ber (0.5)

Arm 2 : Ber (0.8)

$t=1$	$t=2$	$t=3$	$t=4$	$t=5$	(mean for Arm 1 is true no., and mean for Arm 2 is zero)
$I_1 = 1$	$I_2 = 2$	$I_3 = 1$	$I_4 = 1$	$I_5 = 1$	
1	0	0	0		

$\therefore$  Can get stuck with a bad arm

### 2) $\epsilon$ -Greedy Algorithm

Input:  $\epsilon$

$$I_t = \begin{cases} \underset{i}{\operatorname{argmax}} \hat{\mu}_i(t-1) & \text{w.p. } (1-\epsilon_t) \leftarrow \text{exploit} \\ \text{select an arm} \\ \text{uniformly at random} & \text{w.p. } \epsilon_t \leftarrow \text{explore} \end{cases}$$

$\epsilon = 0$  is equivalent to Greedy

$\epsilon = 1$  always select arms uniformly at random.

This will also give linear regret.

Can decrease  $\epsilon$  with increasing  $t$   $\therefore \epsilon_t$  decays as  $t$  increases.

In initial rounds better to explore more, then when  $t$  is large and we have lots of samples of arms then decrease exploration and increase exploitation.

3)

$$P_i = \frac{e^{\eta \hat{\mu}_i(t)}}{\sum_{j=1}^k e^{\eta \hat{\mu}_j(t)}}$$

Similar to weighted majority

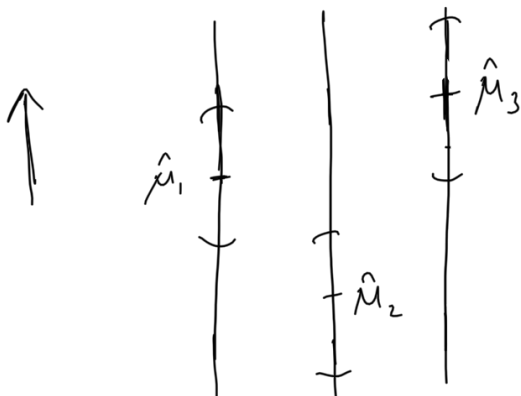
$$P = (P_1, P_2, \dots, P_k)$$

$$I_t \sim P$$

---

## Upper Confidence Bound (UCB) Algorithm

Introduced in 2002



Three arms with 3 estimates of the mean and bounds.

The estimates are obtained after certain number of arms are played.

The larger the bound the lesser number of times that arm has been played.

→ Treat the upper confidence bound as the actual

mean, at any round, and play the arm that has highest upper confidence bound.

→ Thus estimate the confidence intervals and apply greedy algorithm

We know:

$$\hat{\mu}_i = \frac{1}{n} \sum_{t=1}^n x_{ti}$$

$$\Pr \{ \hat{\mu}_i - \mu_i \geq \varepsilon \} \leq \exp \left\{ -\frac{n \varepsilon^2}{2 \sigma^2} \right\}$$

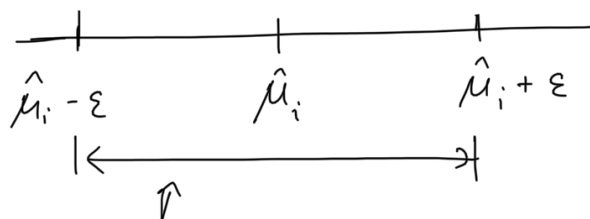
$$\text{Let } \varepsilon = \sqrt{\frac{2 \sigma^2 \log 1/\delta}{n}}$$

Then,

$$\Pr \{ \hat{\mu}_i - \mu_i \geq \varepsilon \} \leq \delta$$

$$\Pr \{ \hat{\mu}_i - \mu_i \leq \varepsilon \} \leq \delta$$

Thus can estimate the confidence bound for  $\mu_i$



true  $\mu_i$  will lie in this interval w.p.  $(1 - 2\delta)$



How does this combine exploration & exploitation:

→ If we have played an arm less then  $T_i(t-1)$  is going to be smaller  $\therefore UCB_i$  will be larger

$\therefore$  may be forced to choose that arm.

But by choosing this arm we are getting a better estimate

$$\begin{array}{ccc} \hat{\mu}_i(t-1) & + & \sqrt{\frac{2 \log t}{T_i(t-1)}} \\ \uparrow & & \uparrow \\ \text{exploit (term i)} & & \text{explore (term ii)} \end{array}$$

→ Due to  $\log t$  term, the exploration never stops, as  $\log t$  keeps increasing, so it makes term ii larger and thus exploration never stops.

→ But  $\log t$  increases slowly, thus sub-optimal arms, even if played will be less no. of times

→ Even when sufficient rounds have been played and a good estimate of  $\hat{\mu}_i$  has been obtained, the exploration never stops due to  $\log t$  term in term ii.