# Stochastic Bandits

Known: number of arms $k$, number of rounds $n > k$

Unknown: $k$ probability distributions $v_1, ..., v_k$ on $[0,1]$

for $t = 1, ..., n$

①  Forcaster chooses $I_t \in \{1, ..., k\}$

②  Given $I_t$, the environment draws reward

$X_{I_t, t} \sim v_{I_t}$ independently from the past

## Regret

$$\max_{i \in [k]} \sum_{t=1}^{n} X_{i,t} - \sum_{t=1}^{n} X_{I_t, t} = R_n$$

$\mu^* = \max_{i \in [k]} \mu_i$

$i^* = \arg\max_{i \in [k]} \mu_i$

## Pseudo Regret

$$\max_{i \in [k]} \mathbb{E}\left[ \sum_{t=1}^{n} X_{i,t} - \sum_{t=1}^{n} X_{I_t, t} \right] = \bar{R}_n \leq \{ \quad \}$$

$$\Rightarrow \bar{R}_n = \mu^* n - \sum_{t=1}^{n} \mathbb{E}[\mu_{I_t}]$$

## UCB1 Strategy

Choose the arm w/ the largest

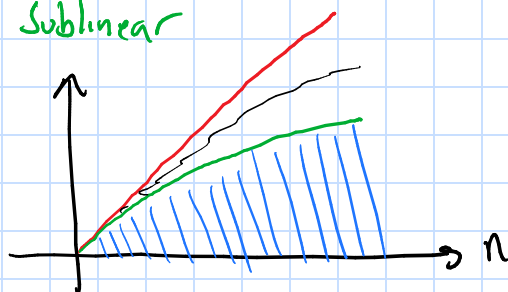$$a_i = \bar{X}_i + \sqrt{\frac{2 \log(t)}{n_i}}$$

$\bar{X}_i \rightarrow$ sample mean of $i^{th}$ arm

$n_i \rightarrow$ # of times $i^{th}$ arm was sampled

**Thm:** Let $\Delta_i = \mu^* - \mu_i$. For all $k > 1$, if UCB1 is run on $k$ machines having arbitrary distributions $v_1, \ldots, v_k$ then its expected regret after $n$ plays is at most

$$\left[ 8 \sum_{i: \mu_i \neq \mu^*} \frac{\log(n)}{\Delta_i} \right] + \left(1 + \frac{\pi^2}{3}\right) \sum_{j=1}^{k} \Delta_j$$

☆ Sublinear



| | Br/round | cumulative |
|---|---|---|
| $t_1$ | $r_1$ | $r_1$ |
| $t_2$ | $r_2$ | $r_1 + r_2$ |
| $t_3$ | $r_3$ | $r_1 + r_2 + r_3$ |
| $t_4$ | $r_4$ | $r_1 + r_2 + r_3 + r_4$ |

**ε - Greedy**    $\varepsilon \rightarrow (0, 1)$

① $S_i \sim Uni(0, 1)$

② if $S_i \geq \varepsilon$

    ☆ Choose $I_t$ to be the largest mean

  else

    ☆ Choose $I_t$ at random

# The Adversarial Bandit

**known parameters:** number of arms $k \geq 2$, # of rounds

for $t = 1, \ldots, n$

    ① The forecaster chooses $I_t \in \{1, \ldots, k\}$

    ② The adversary chooses a gain $g_t = (g_{1t}, g_{2t} \ldots g_{kt})$

    ③ The forecaster gets reward $g_{I_t t} \in [0,1]^k$

       (none of the other gains are revealed)

# Exp3

Inputs: $\gamma \in [0,1]$, $W_i(1) = 1$    $i \in [k]$, $n$

for $t = 1, \ldots n$

    ① $P_i(t) = (1-\gamma) \dfrac{W_i(t)}{\sum_j W_j(t)} + \dfrac{\gamma}{k}$

    ② Draw $I_t$ from $P_1(t), \ldots, P_k(t)$

    ③ $X_{I_t, t} \in [0, 1]$

    ④ for $j = 1, \ldots, k$

    ⑤ $\hat{X}_j(t) = \begin{cases} X_j(t)/P_i(t) & j = I_t \\ 0 & \text{else} \end{cases}$

    ⓐ $W_j(t+1) = W_j(t) \exp\left(\gamma \hat{X}_j(t) / k\right)$