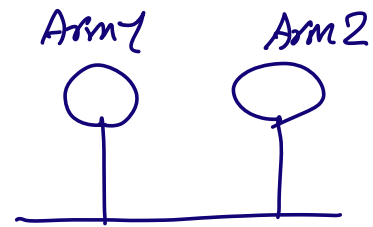Lecture 2 - Jan 17, 2023

- Multi-Arm Bandits

  - Introduction
  - Exploration - Exploitation Delima
  - Epsilon - Greedy Policy
  - Optimistic Initial Values

  - Upper Confidence Bound Selection Policy
  - Gradient - Based Selection Policy
  - Thompson Sampling

HW1 is assigned ⟶ Due Jan 27
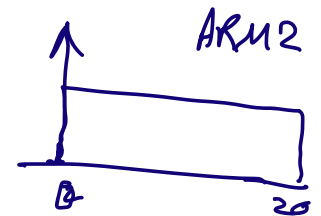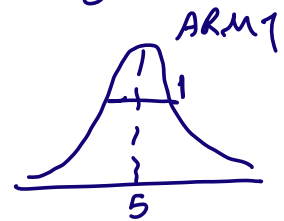
TA's first office hour: Friday, Jan 20, 12pm - 1pm

# Multi-Arm Bandit

Arm 1    Arm 2

$ARM\underline{1} = Reward \quad R^1 | a^1 \sim N(\mu=5, \sigma=1)$

$ARM2 = \qquad R^2 | a^2 \sim Uniform \begin{bmatrix} a & b \\ 0 & 20 \end{bmatrix}$

ARM1

$\Rightarrow$ Goal: Maximizing the total reward

5

ARM2

$\overset{*}{Q}(a) = E[R_t | a_t = a] \qquad a \in \mathcal{A} = \{a^1, a^2, \cdots, a^k\}$

## Averaging Learning Rule:

$Q_2(a) = R_1$

$Q_3(a) = \dfrac{R_1 + R_2}{2}$

$$Q_n^{(a)} = \frac{R_1 + R_2 + \cdots + R_{n-1}}{n-1}$$

$$Q_{n+1}^{(a)} = \frac{R_1 + R_2 + \cdots + R_{n-1} + R_n}{n} = \frac{R_1 + R_2 + \cdots + R_{n-1}}{n} + \frac{R_n}{n} \times \frac{n-1}{n-1}$$

$$\times \frac{n-1}{n-1}$$

$$= Q_n(a)\,\frac{n-1}{n} + \frac{1}{n}\,R_n$$

Last reward

$$Q_{n+1}(a) = Q_n(a) + \frac{1}{n}\left[\overbrace{R_n} - Q_n(a)\right] \qquad (\mathrm{I})$$

Previous estimate

Previous estimate

$n = 10 \longrightarrow$

$$Q_n(a) + \frac{1}{10}\left[\bigcirc - Q_n(a)\right]$$

$n = 1000\,000 \longrightarrow$

$$Q_n(a) + \frac{1}{1000\,000}\left[\bigcirc - Q_n(a)\right]$$

# Nonstationary

$$Q(a) \leftarrow Q(a) + \alpha \left[ R - Q(a) \right]$$

Learning Rate

$$\sum_{t=1}^{\infty} \alpha_t = \infty \qquad \sum_{t=1}^{\infty} \alpha_t^2 < \infty \implies Q^*(a)$$

$$\alpha_t = \frac{1}{t}$$

$\alpha = 0.1$

$R = 100 \qquad Q(a) = 0$

## Policy 1: Epsilon-Greedy Policy
### $\varepsilon$-greedy



$$Q(a^1) = Q(a^2) = \cdots = Q(a^k) = 0$$

$$a_t \sim \begin{cases} \text{Greedy.} \quad \underset{a \in \{a^1, \dots, a^k\}}{\arg\max} \, Q(a) & \text{w.p. } 1-\varepsilon \implies \text{Exploitation} \\[2mm] \text{Random} \{a^1, a^2, \dots, a^k\} & \text{w.p. } \varepsilon \implies \text{Exploration} \end{cases}$$

**Example:** $\varepsilon = 0.1$
$\alpha = 0.5$

$R \sim N(5, \sigma = 10)$

$R^2 \sim Unif[8, 12]$

$Q(a^1) = Q(a^2) = 0$

①

$\pi^{\varepsilon\text{-greedy}} = \begin{cases} \underset{a \in \{a^1, a^2\}}{argmax} \; Q(a) & w.p. \; 0.9 \\ \\ Random \; \{a^1, a^2\} & w.p. \; 0.1 \end{cases} \longrightarrow a^1 \rightsquigarrow$

$argmax \{ \overset{a^1}{\underset{-}{0}}, \; \overset{a^2}{\underset{-}{0}} \}$

Random    greedy

$a^1 \longleftarrow 0.5 \longrightarrow a^2$

$R = 10 \sim N(5, 10)$

$Q(a^1) = Q(a^1) + \alpha [R - Q(a^1)]$

$= \quad 0 \quad + 0.5 [10 - 0] = 5$

$$\boxed{Q(a^1)=5, \quad Q(a^2)=0}$$

② 

$\pi^{\text{E-greedy}} = \begin{cases} \text{argmax} \left\{ \overset{a^1}{\overline{5}}, \overset{a^2}{\overline{0}} \right\} = a^1 \text{ w.p. 0.9} \\ \\ \text{Random} \qquad\qquad\qquad 0.1 \end{cases} \sim a^1$

$R = 3 \curvearrowleft N(5, 10)$

$Q(a^1) = \underline{Q(a^1)} + \alpha \left[ R - Q(a^1) \right]$

$\qquad = 5 + 0.5 \left[ 3 - 5 \right] = 4$

$$\boxed{Q(a^1)=4, \quad Q(a^2)=0}$$

③

$\pi^{\text{E-greedy}} = \begin{cases} \text{argmax } Q(a) = a^1 \qquad 0.9 \\ \\ \text{Random} \{a^1, a^2\} \qquad 0.1 \end{cases} \sim \underset{\text{Exploration}}{a^2}$

$R = 9 \sim \text{Uniform } [8, 12]$

$$Q(a^2) = Q(a^2) + \alpha [\underline{R} - Q(a^2)]$$

$$= 0 + 0.5 [9 - 0] = 4.5$$

$$\boxed{Q(a^1) = 4, \quad Q(a^2) = 4.5}$$

④

$$\pi^{\varepsilon\text{-}g} = \begin{cases} \text{argmax } Q(a) = a^2 & 0.9 \\ \text{Radm}\{a^1, a^2\} & 0.1 \end{cases}$$

Reward

$Q^*(a^4)$

$a^1$  $a^2$  $a^3$  $a^4$

Ave Accumulated Reward

$\varepsilon$

$0$   $1$

$\begin{cases} a^1 \\ \vdots \\ a^K \end{cases}$

$a^i$  $a^i$ ----- $a^i$ $a^b$ $a^?$ -------

$a_t \sim \begin{cases} \text{greedy. } \arg\max\limits_{a \in \{a^1,...,a^K\}} Q(a) & \text{w.p. } 1-\varepsilon \Rightarrow \text{Exploitation} \\ \\ \text{Random} \{a^1, a^2, ..., a^K\} & \text{w.p. } \varepsilon \Rightarrow \text{Exploration} \end{cases}$

Reward

$\varepsilon = 0.1$

$\varepsilon = 0.3$

$\varepsilon = 0.02$

$r_2$

$r_1$

Time / step

episode 1

$a^1$    $a^1$

$r_1$    $r_2$