# Lecture 3 - Jan 20, 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 → Due Jan 27

Project 1 is Posted → Due Feb 7

TA's first office hour: Friday, Jan 20, 12pm-1pm

overview:

$$Q^*(a) = E[R \mid a]$$

$$a^* = \arg\max_{a \in A} Q^*(a)$$

Learning ⟵ Distribution are unknown

⟶ Policy: $\underline{a}$

⟶ Estimation $Q(a)$

$Q(a) = 0$ for all $a \in A$

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

$a^1 \quad a^2 \quad a^k$

$R \mid a$ ⟵ Random variable

$a^1$

$R \mid a^1$

$a^1 \quad a^2 \quad a^k$

Policy 1: $\varepsilon$-greedy

$$a \sim \begin{cases} \arg\max_{a \in A} Q(a) \\ \text{Random}\{a^1, \ldots, a^k\} \end{cases}$$

$\sigma \quad 5 \quad 100$

$1 - \varepsilon$ ⟵ greedy

$\varepsilon$ ⟵ Random

$a^1 \longrightarrow \mathbb{Q}(a^1) = 20, \ Q(a^2) = 0$

$\downarrow$ $a^2 \xrightarrow{R=90} Q(a^1) = 20, \ Q(a^2) = 45$

$40 \curvearrowleft$



---

**Optimistic Initial Value:** $Q^*(a^1) = 5$
$Q^*(a^2) = 6$

$\mathbb{Q}(a^1) = \mathbb{Q}(a^2) = 0$



Time 1

action $a^1$

Reward 5

$$Q(a^1) = \underset{0}{\underline{Q(a^1)}} + \overset{0.5}{\alpha}\,[\underset{5}{R} - \underset{0}{\underline{Q(a^1)}}] = 2.5$$

$\pi^{\varepsilon-g} \longrightarrow \begin{cases} \text{greedy} = a^1 \quad 1-\varepsilon \\ \text{Random} \{a^1, a^2\} \ \varepsilon \end{cases}$

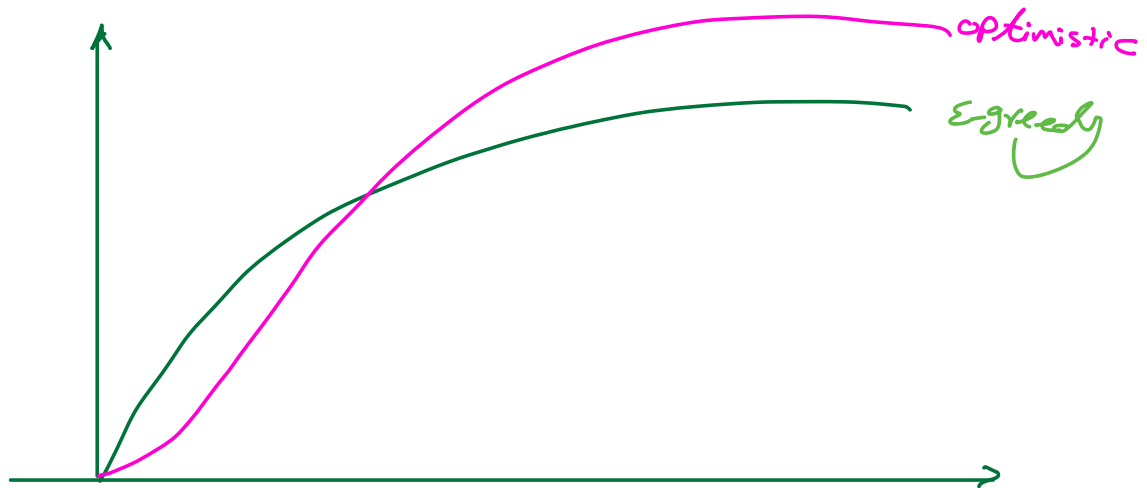$\varepsilon = 0.1 \longrightarrow \begin{cases} a^1 \to 0.95 \\ a^2 \to 0.05 \end{cases}$

$\boxed{\begin{array}{l} Q(a^1) = 2.5 \\ Q(a^2) = 0 \end{array}}$

$\begin{cases} Q(a^1) = Q(a^2) = 15 \\[4pt] \underset{15}{Q(a^1)} = \underset{15}{Q(a^1)} + \underset{0.5}{\alpha}\,[\underset{5}{R} - \underset{15}{Q(a^1)}] = 10 \\[6pt] \pi^{\varepsilon-g} \to \begin{cases} \text{greedy} = a^2 \\ \text{Random} \{a^1, a^2\} \end{cases} \quad \begin{array}{l} Q(a^1) = 10 \\ Q(a^2) = 15 \end{array} \\[10pt] \varepsilon = 0.1 \begin{cases} a^1 \leftarrow 0.05 \\ a^2 \leftarrow 0.95 \end{cases} \end{cases}$
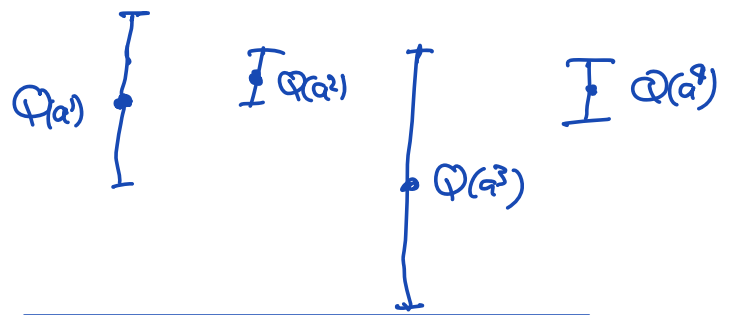
optimistic

ε-greedy

---

# Upper Confidence Bound (UCB) ⇐ Policy

$$a_t = \underset{a \in \{a', a^2, \dots, a^k\}}{\arg\max} \left[ Q(a) + c \sqrt{\frac{\log t}{N_t(a)}} \right]$$

addition to the current estimate of $Q(a)$

$t$: time step

$N_t(a)$: # that action $a$ is selected up to time $t$

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

$$Q(a^1) = 10 \qquad\qquad \widehat{Q}(a^2) = 1$$
$$N(a^1) = 100 \qquad\qquad N(a^2) \leq 1 \quad \leadsto \quad t = 101$$

$$t = 102 \longrightarrow a_t = \arg\max \left[ Q(a) + c \sqrt{\frac{\log t}{N_t(a)}} \right]$$

$$\underset{a \in \{a^1, a^2\}}{\arg\max} \left\{ \underbrace{10 + c \sqrt{\frac{\log 102}{100}}}_{a^1} \quad , \quad \underbrace{1 + c \sqrt{\frac{\log 102}{1}}}_{a^2} \right\}$$

$$c = 1 \leftarrow \quad 10 + 0.213 \qquad\qquad 1 + 1.52$$

$$c = 10 \leftarrow \quad 10 + 2.13 \qquad\qquad 1 + 15.2$$

$$c = 0 \implies \text{greedy}$$

c larger $\longrightarrow$ More exploration

Online Learning Algorithms $\begin{cases} \text{Exp3} \\ \text{Hedge} \\ \text{Regret Matching} \end{cases}$

# Policy: Gradient-Bandit Policy

Directly parametrizing the policy:

$H_t(a)$: numeric preference for action $a$

$$P(a_t = a) = \frac{e^{H_t(a)}}{\sum_{c=1}^{K} e^{H_t(a^c)}} : \pi_t(a)$$

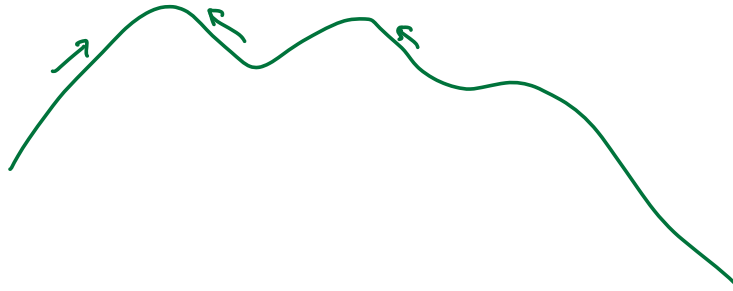Probability of taking action $a$ at time $t$

Boltzman or Softmax Policy

Gibbs distribution

$$\left.\begin{array}{l} H_1(a^1) = 0 \\ H_1(a^2) = 0 \end{array}\right\} \rightarrow \left\{\begin{array}{l} \pi_1(a^1) = \dfrac{e^{H_1(a^1) = 0}}{e^{H_1(a^1)} + e^{H_1(a^2)}} = \dfrac{1}{1+1} = \dfrac{1}{2} \\[2em] \pi_1(a^2) = \dfrac{e^0}{e^0 + e^0} = \dfrac{1}{2} \end{array}\right.$$

$$\left.\begin{array}{l} H_5(a^1) = 8 \\ H_5(a^2) = 0 \end{array}\right\} \rightarrow \left\{\begin{array}{l} \pi_5(a^1) = \dfrac{e^8}{e^8 + e^0} = 0.9997 \\[2em] \pi_5(a^2) = \dfrac{e^0}{e^8 + e^0} = 0.0003 \end{array}\right.$$

$$H_{t+1}(a) = H_t(a) + \alpha \frac{\partial E[R_t]}{\partial H_t(a)}$$

← Average reward



$$\begin{cases} H_{t+1}(A_t) = H_t(A_t) + \alpha (R_t - \bar{R}_t)(1 - \pi_t(A_t)) \\ H_{t+1}(a) = H_t(a) - \alpha (R_t - \bar{R}_t)\pi_t(a) \quad \text{for all } a \in A - A_t \end{cases}$$

$A_t$: action selected at time $t$

$a$: all actions expect $A_t$

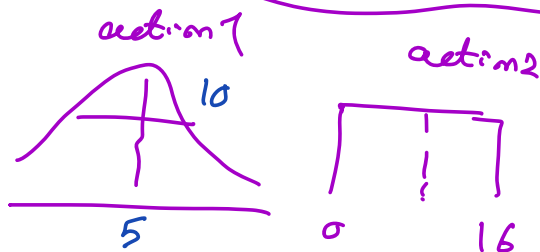$\bar{R}_t$ : Average Reward up to current time $t$ including "$t$"

$\bar{R}_t = 0$ , $R_t = 10$ , $A_t = a^1$ $\qquad A = \{a^1, a^2\}$

Action 1 → $\quad + \alpha (\underbrace{R_t}_{10} - \underbrace{\bar{R}_t}_{0}) (1 - \underbrace{\pi_t(A_t)}_{0.6})$ $\qquad \pi_t(a^1) = 0.6$

$\pi_t(a^2) = 0.4$

action 2 → $\quad - \alpha (\underbrace{R_t}_{10} - \underbrace{\bar{R}_t}_{0}) \underbrace{\pi_t(a)}_{0.4}$

Action 1 → $\quad + \alpha \, (\underbrace{R_t}_{10} - \underbrace{\bar{R}_t}_{0})(1 - \underbrace{\pi_t(A_t)}_{0.9999})$ $\quad \overbrace{\dfrac{e^{H_t(a)}}{\sum_a e^{H_t(a^c)}}}^{0.0001}$

$\pi_t(a') = 0.9999$
$\pi_2(a') = 0.0001$

action 2 → $\quad -\alpha \, (\underbrace{R_t}_{10} - \underbrace{\bar{R}_t}_{0})\underbrace{\pi_t(a)}_{0.0001}$

---

action 1



5

action 2

0       16

$R^{a'} \sim \mathcal{N}(5, 10)$

① 

$H_1(a') = H_1(a^2) = 0$

↓ policy

$\pi_1(a') = \dfrac{e^{H_1(a')}}{e^{H_1(a')} + e^{H_1(a^2)}} = \dfrac{e^o}{e^o + e^o} = \dfrac{1}{2}$

$\pi_1(a^2) = \dfrac{e^o}{e^o + e^o} = \dfrac{1}{2}$

$\rightsquigarrow A_1 = a'$

$A_1 = a', \ R_1 = 10 \quad \Rightarrow \bar{R}_1 = \dfrac{\boxed{10}}{1} = 10$ $\qquad \overbrace{\dfrac{r_1 + r_2 + r_3 + \cdots + r_T}{T}}$

preference update

$H_2(a') = \underbrace{H_1(a')}_{0} + \underbrace{\alpha}_{0.5} [\underbrace{R_1}_{10} - \underbrace{\bar{R}_1}_{10}](1 - \underbrace{\pi_1(a')}_{1/2}) = 0$

$H_2(a^2) = \underbrace{H_1(a^2)}_{0} - \underbrace{\alpha}_{0.5}[\underbrace{R_1}_{10} - \underbrace{\bar{R}_1}_{10}]\underbrace{\pi_1(a^2)}_{1/2} = 0$

② 

$H_2(a^1) = 0$ , $H_2(a^2) = 0$

$\quad\quad \hookrightarrow \pi_2(a^1) = \dfrac{e^0}{e^0 + e^0} = \dfrac{1}{2}$ $\left.\phantom{\dfrac{1}{2}}\right\}$ $\longrightarrow A_t = a^2$

$\quad\quad\quad \pi_2(a^2) = \dfrac{1}{2}$

$A_2 = a^2,\ R_2 = 2 \quad\quad \Rightarrow \bar{R}_2 = \dfrac{10 + 2}{2} = 6$

preloace Upalate

$\left\{\begin{array}{l}
H_3(a^2) = \underbrace{H_2(a^2)}_{0} + \underbrace{\alpha}_{0.5}\left[\underbrace{R_2}_{2} - \underbrace{\bar{R}_2}_{6}\right]\left(1 - \underbrace{\pi_2(a^2)}_{1/2}\right) = -1 \\[20pt]
H_3(a^1) = \underbrace{H_2(a^1)}_{0} - \underbrace{\alpha}_{0.5}\left[\underbrace{R_2}_{2} - \underbrace{\bar{R}_2}_{6}\right]\underbrace{\pi_2(a^1)}_{1/2} = 1
\end{array}\right.$

_____

$H_3(a^1) = 1$ , $H_3(a^2) = -1$

$\quad \hookrightarrow \pi_3(a^1) = \dfrac{e^1}{e^1 + e^{-1}} = 0.88$ $\left.\phantom{\dfrac{1}{2}}\right\}$

$\quad\quad\quad \pi_3(a^2) = \dfrac{e^{-1}}{e^1 + e^{-1}} = 0.12$