Lecture 4 - Jan 24, 2023

- Multi Arm Bandits
  - Introduction
  - Exploration - Exploitation Delima
  - Epsilon-Greedy Policy
  - Optimistic Initial Values
  - Upper Confidence Bound Selection Policy
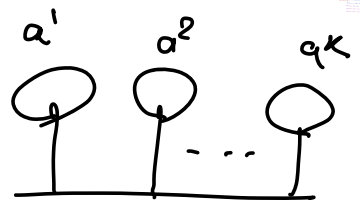  - Gradient-Based Selection Policy
  - Thompson Sampling
- Reinforcememt Learning Preliminaries
  - State, Action, Reward, Policy

HW1 → Due Jan 28
Project 1 → Due Feb 7

TA's office hour:   Wendsdays, 12pm-1pm (in-person)
                    Fridays, 12pm-1pm (virtual)

Policy          Q-Estimate



a¹   a²   ... aᵏ

1- $\varepsilon$-greedy

$$\mathcal{L} \quad a \sim \begin{cases} \underset{a \in A}{\operatorname{argmax}} \; Q(a) & \text{w.p } 1-\varepsilon \\ \\ \text{Random } \{a^1, ..., a^k\} & \varepsilon \end{cases}$$

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

$\rightarrow$ Larger $\varepsilon \rightarrow$ more exploration

Optimistic initial Value $\Rightarrow$ $Q(a) = 0$ $\nearrow$ large Value

2- Upper Confidence Bound Policy

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

$$a_{t+1} = \underset{a \in A}{\operatorname{argmax}} \left[ Q(a) + c \sqrt{\frac{\log t}{N_a(t)}} \right]$$

Larger $c \longrightarrow$ More exploration
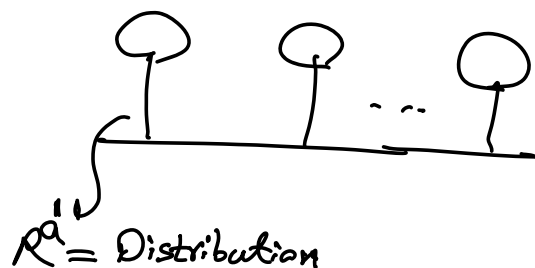
3- Gradient-Bandit policy ( does not $Q$ )

$\rightarrow$ preference
$$H_0(a) = 0 \quad \text{for all } a \in A$$

$$\pi_t(a) = \frac{e^{H_t(a)}}{\sum_{b \in A} e^{H_t(b)}}$$

$$\begin{cases} H_{t+1}(A_t) = H_t(A_t) + \alpha [R_t - \bar{R}_t](1 - \pi_t(A_t)) \\ \\ H_{t+1}(a) = H_t(a) - \alpha [R_t - \bar{R}_t] \pi_t(a) \end{cases}$$

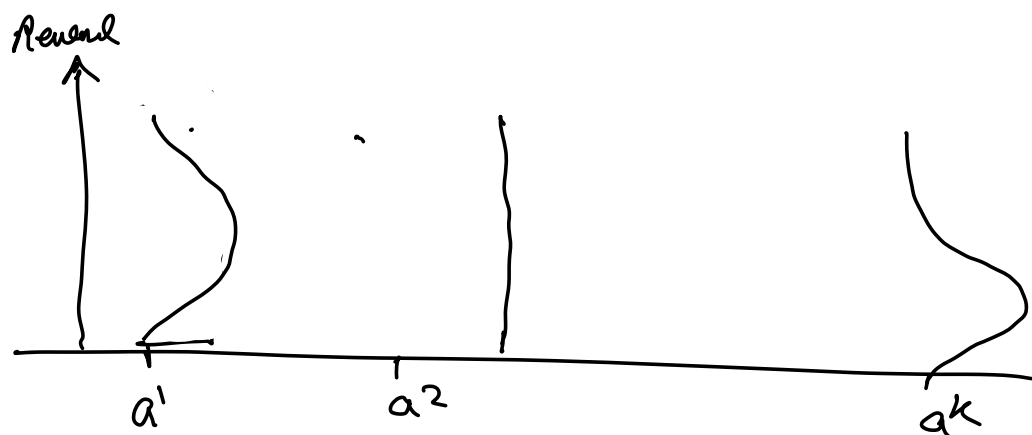# 4- Thompson Sampling → Bayesian Bandit Algorithms

---

Bernoulli (P)
  0.9

$R^{a^i}$:

$R_i^{a^i}$: Bernoulli $\left( \dfrac{\alpha_i}{\alpha_i + \beta_i} \right)$

$\underbrace{\qquad}$

$R^{a^i}$ = Distribution

Beta Distribution for likelihood of observed win/loss

Beta$( \alpha_i, \beta_i )$

arem 2 selected

$\dfrac{Win}{loss} \rightarrow \alpha_2 = \alpha_{2}' + 1$

loss $\rightarrow \beta_2 = \beta_2 + 1$

$P(\theta_a \mid D) = \dfrac{P(D \mid \theta^a) \, P(\theta^a)}{P(D)}$

Reward

$a^1 \qquad a^2 \qquad a^k$

$$r_1 \sim R^{a^1} \qquad r_2 \sim R^{a^2} \qquad r_k \sim R^{a^k}$$

$$a_k = \underset{c}{\text{argmax}} \; r_i$$

---

# Proof of Gradient Bandit Policy

$$\pi_t(a) = \frac{e^{H_t(a)}}{\sum_{b \in A} e^{H_t(b)}}$$



$$H_t(a^1) = 0 \atop H_t(a^2) = 0 \; \rightarrow \; \begin{cases} \pi_t(a^1) = \frac{1}{2} \\ \pi_t(a^2) = \frac{1}{2} \end{cases}$$

$$H_t(a^1) = 8 \atop H_t(a^2) = 0 \; \rightarrow \; \begin{cases} \pi_t(a^1) = 0.9997 \\ \pi_t(a^2) = 0.0003 \end{cases}$$

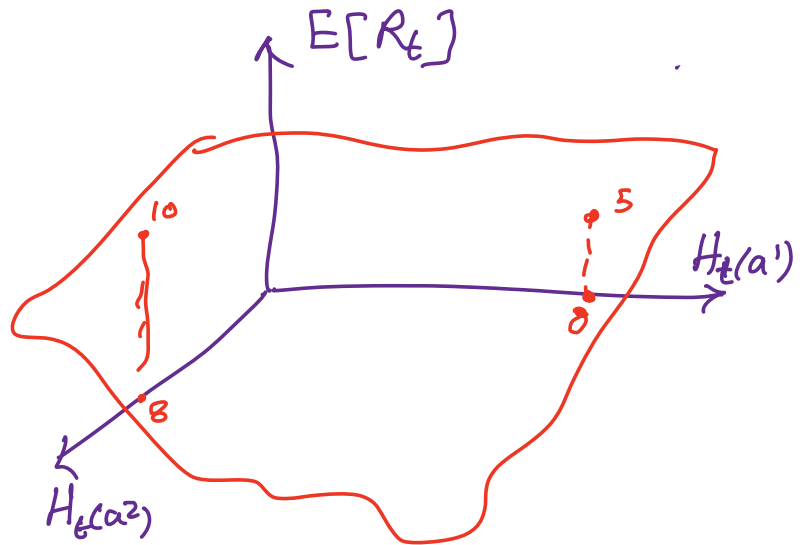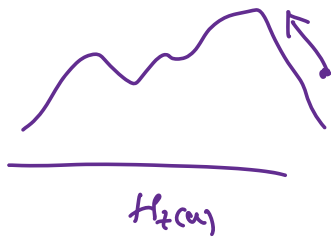Objective : maximizing Accumulated
Reward

$$E[R_t]$$
$$\searrow \; E[R_t \mid \pi_t(a^1), \pi_t(a^2)]$$

$$= E[R_t \mid H_t(a^1), H_t(a^2)]$$

$$Q(a^1) = 5$$
$$Q(a^2) = 10$$

$$= \sum_{b \in A} Q(b) \, \pi_t(b)$$

$$H_{t+1}(a) = H_t(a) + \alpha \frac{\partial E[R_t]}{\partial H_t(a)} \quad \text{for all } a \in A$$

$$\frac{\partial E[R_t]}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \left[ \sum_{b \in A} Q(b) \, \pi_t(b) \right]$$

$$= \sum_{b \in A} Q(b) \frac{\partial \pi_t(b)}{\partial H_t(a)}$$

$$= \sum_{b \in A} \left( Q(b) - \overset{R_t}{X_t} \right) \frac{\partial \pi_t(b)}{\partial H_t(a)}$$

any scalar

why?

$$\pi_t(a^1) = \frac{e^{H_t(a^1)} x}{e^{H_t(a^1)} + e^{H_t(a^2)} y}$$

$$\pi_t(a^2) = \frac{e^{H_t(a^2)}}{e^{H_t(a^1)} + e^{H_t(a^2)}}$$

$$\rightarrow \sum_{b \in A} \frac{\partial \pi_t(b)}{\partial H_t(a^1)} = \frac{\partial}{\partial x} \frac{e^x}{e^x + e^y} + \frac{\partial}{\partial x} \frac{e^y}{e^x + e^y}$$

$$\frac{\partial E[R_t]}{\partial H_t(a)} = \sum_{b \in A} \pi_t(b) \left( Q(b) - \bar{R}_t \right) \frac{\partial \pi_t(b)}{\partial H_t(a)} \Big/ \pi_t(b)$$

$$= E\left[ \left( Q(b) - \bar{R}_t \right) \frac{\partial \pi_t(b)}{\partial H_t(a)} \Big/ \pi_t(b) \right]$$

$$A_t \sim \pi_t = \pi_t(a^1), \pi_t(a^2)$$

$$\frac{\partial E[R_t]}{\partial H_t(a)} \simeq \underbrace{\left( Q(A_t) - \bar{R}_t \right)}_{R_t} \underbrace{\frac{\partial \pi_t(A_t)}{\partial H_t(a)}}_{} \Big/ \pi_t(A_t)$$

$$\frac{\partial \pi_t(A_t)}{\partial H_t(a)} \nearrow \quad a \neq A_t \cdot \quad \pi(A_t)\left( - \pi(a) \right)$$

$$\searrow \quad a = A_t \quad \rightarrow \quad \pi(A_t)\left( 1 - \pi(A_t) \right)$$

$$\pi(A_t) = \frac{e^x}{e^x + e^y}$$

$$\pi(a) = \frac{e^y}{e^x + e^y}$$

$$\longrightarrow \frac{\partial \pi_t(A_t)}{\partial H_t(A_t)} = \frac{\partial}{\underbrace{\partial H_t(A_t)}_{x}} \frac{e^x}{e^x + e^y}$$

$$\begin{cases} H_{t+1}(A_t) = H_t(A_t) + \alpha\,[R_t - \bar{R}_t]\,(1 - \pi_t(A_t)) \\[2mm] H_{t+1}(a) = H_t(a) - \alpha\,[R_t - \bar{R}_t]\,\pi_t(a) \end{cases}$$

# Reinforcement Learning - preliminaries

## State space $S$   $s \to$ state

$s \in S$

| 9 | 10 | 11 | 12 |
|---|----|----|----|
| 8 | ■ | 14 | 13 |
| 7 | ■ | 16 | 15 |
| 6 | 5 | ■ | ■ |
| 4 | 3 | 2 | 1 |

■ Wall  ▢ Bump  ▢ Goal

## Action Space $A$

$a \in A$

$A = \{ UP, Down, Left, Right\}$

$$R: S \times A \times S \to Real$$

Immediate Reward $R(s, a, s')$

$$\begin{cases} -1 & \text{any movement} \\ -10 & \text{Bump} \\ 100 & \text{Goal} \end{cases}$$

$R(3, a = L, 4) = -1 - 10 = -11$

$\underbrace{}_{movement} \underbrace{}_{Bump}$

$R(3, U, (1))$

Goal: Find the best $\overset{\text{sequence of actions}}{\text{path}}$ that results in highest Accumulated Reward

| 9 | 10 | 11 | 12 |
|---|----|----|----|
| 8 | ■ | 14 | 13 |
| 7 | ■ | 16 | 15 |
| 6 | 5 | ■ | ■ |
| 4 | 3 | 2 | 1 |

■ Wall  ▢ Bump  ▢ Goal

$$\underset{a_{0:T}}{\text{argmax}} \; E\left[ \sum_{t=0}^{T-1} R(s_t, a_t, s_{t+1}) \mid s_0 = 2, \; a_{0:T-1}, \; s_{t+1} \sim P(s' \mid s_t, a) \right]$$

$2 \to \bigcirc \to \bigcirc$
$\quad a_0 \quad - a_1 \; - -$

$a_{20}$ $\overset{ST}{\bigcirc}$

| 9 | 10 | 11 | 12 |
|---|----|----|----|
| 8 | ⬛ | 14 | 13 |
| 7 | ⬛ | 16 | 15 |
| 6 | 5 | ⬛ | ⬛ |
| 4 | 3 | 2 | 1 |

⬛ Wall   🟧 Bump   🟩 Goal

-1 -1 -1 -1 -1 -1 -- +99 +100 +100 —

$\underbrace{\phantom{-1 -1 -1 -1 -1 -1 -- }}$
-11

+99 +

$\underbrace{\phantom{+99 + }}$
800

-20

| 9 | 10 | 11 | 12 |
|---|----|----|----|
| 8 | ⬛ | 14 | 13 |
| 7 | ⬛ | 16 | 15 |
| 6 | 5 | ⬛ | ⬛ |
| 4 | 3 | 2 | 1 |

⬛ Wall   🟧 Bump   🟩 Goal