

8/8/19

## AI5001: Intro to Modern AI

- Review & discussion ( $A_{t+1} \rightarrow A_t$ )
- Exploration vs exploitation:  $\epsilon$ -greedy method
- Estimating action value: simple averaging, incremental, non stationarity
- $A_t = \underset{a}{\operatorname{argmax}} Q_t(a)$  : greedy
- $\epsilon$ -greedy method:  $(1 - \epsilon)$  times for exploitation,  
 $\epsilon$  times for exploration ( $0 < \epsilon < 1$ )

→ HW problem to understand this better.

### Estimating action value:

$a$ : actions

$q(a)$ : true value for 'a'

$Q_t(a)$ : estimated value at 'a' after 't' timesteps

Assume: stochastic nature of reward. (Stationary or rewards coming from a fixed distribution)

Recall: Our goal is to maximise expected value/reward

$$Q_t(a) = \frac{1}{N_t(a)} \cdot [R_1(a) + R_2(a) + \dots + R_{N_t(a)}(a)]$$

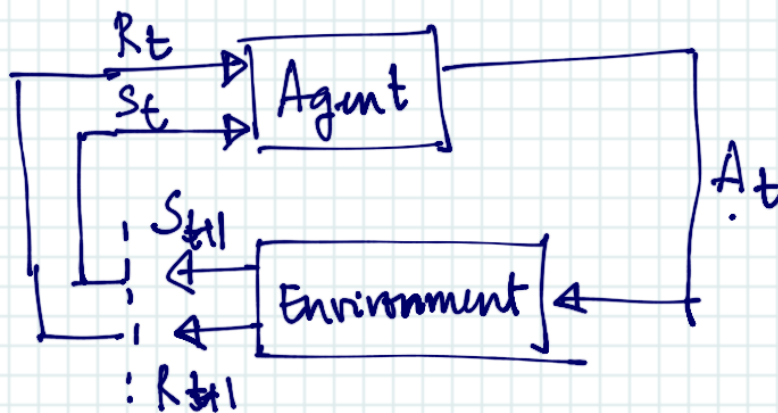
Note: this operation is expensive in terms of memory, so use incremental method

$$\begin{aligned}
 Q_t &= \frac{1}{N_t} \sum_{i=1}^{N_t} R_i \\
 &= \frac{1}{N_t} \left[ R_{N_t} + \sum_{i=1}^{N_t-1} R_i \right] \\
 &= \frac{R_{N_t}}{N_t} + \frac{N_t-1}{N_t} \cdot Q
 \end{aligned}$$

$$Q_t = Q_{t-1} + \frac{1}{N_t} [R_{N_t} - Q_{t-1}]$$

Nonstationary bandits: Weighted sum of rewards. We weight nearer time rewards (recent rewards) more than older rewards

Define the RL problem:



Policy:  $\pi_t(a/s)$