# Lecture 19 - March 28, 2023

- Temporal Difference Learning
  - SARSA Q-Learning
  - On-Policy Vs. Off-Policy
  - Expected SARSA Double Q-Learning
  - Multi-step Bootstrapping
  - SARSA-Lambda
  - Actor-Critic Method
- Function Approximation in Reinforcement Learning
  - Basics of Function Approximations
  - Least Square Policy Iteration (LSPI)
  - Neural Fitted Q-Iterations (NFQI)

---
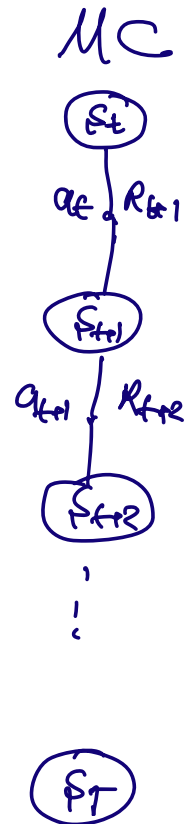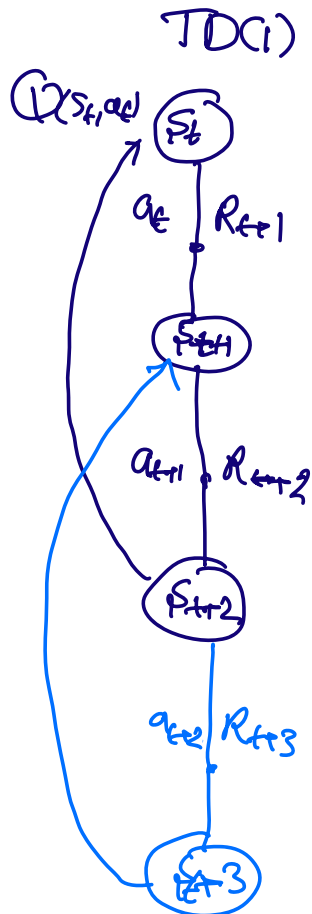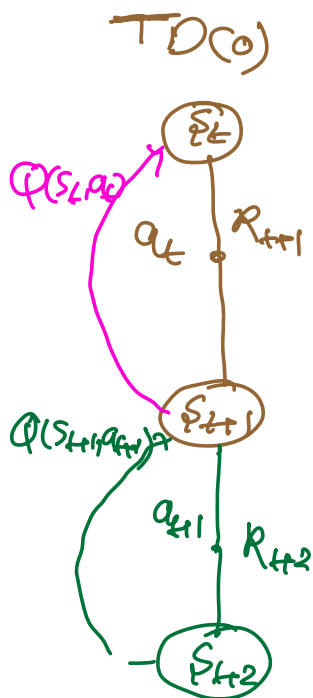
HW 4 → Due March 31

Exam 2 → Tues, April 4

Project 3 → Due April 14

TA's office hour:　Wendsdays, 2pm-3pm (in-person)
　　　　　　　　　Fridays, 2pm-3pm (virtual)

# Multi-Step Boot strapping

## TD(0)

$Q(S_t, a_t)$

$Q(S_{t+1}, a_{t+1})$

$S_t$

$a_t$    $R_{t+1}$

$S_{t+1}$

$a_{t+1}$    $R_{t+2}$

$S_{t+2}$

## TD(1)

$Q(S_t, a_t)$

$S_t$

$a_t$    $R_{t+1}$

$S_{t+1}$

$a_{t+1}$    $R_{t+2}$

$S_{t+2}$

$a_{t+2}$    $R_{t+3}$

$S_{t+3}$

## MC

$S_t$

$a_t$    $R_{t+1}$

$S_{t+1}$

$a_{t+1}$    $R_{t+2}$

$S_{t+2}$

$\vdots$

$S_T$

TD(0):

$$Q(S_t, a_t) = Q(S_t, a_t) + \alpha \left[ R_{t+1} + \gamma \max Q(S_{t+1}, a) - Q(S_t, a_t) \right]$$

TD(1)

$$Q(S_t, a_t) = Q(S_t, a_t) + \alpha \left[ R_{t+1} + \gamma R_{t+2} + \gamma^2 \max_a Q(S_{t+2}, a) - Q(S_t, a_t) \right]$$

TD(n) longer variance $\Longleftarrow$ Variance
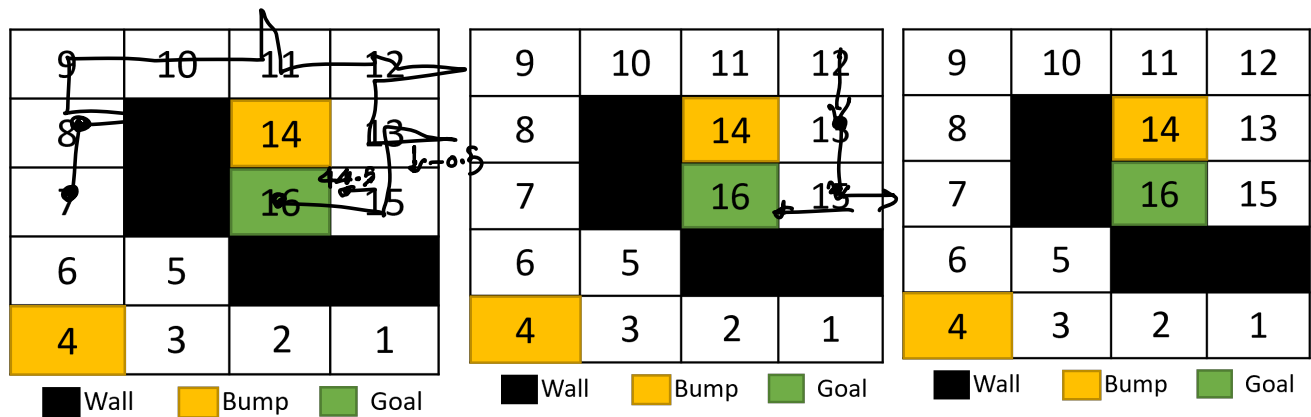
TD(n) Less Biased $\Longleftarrow$ Biased
where n>0

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

⬛ Wall   🟧 Bump   🟩 Goal

TD(1)

$$Q(13, D) = Q(13, D) + \alpha \left[ -\overset{R_{tot}}{T} + \gamma(99) + \gamma^2 Q(16, a_i) - Q(13, D) \right]$$

# SARSA-Lambda — SARSA-$\lambda$



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

■ Wall  🟧 Bump  🟩 Goal

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

■ Wall  🟧 Bump  🟩 Goal

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

■ Wall  🟧 Bump  🟩 Goal

$$Q(7,U) = \underbrace{Q(7,U)}_{} + \underbrace{\alpha}_{0.5} \Big[ -1 + \gamma \underbrace{Q(8,R)}_{} - \underbrace{Q(7,U)}_{} \Big] = -0.5$$

---

## Eligibility Trace $e(s,a)$

Larger $\longrightarrow$ More visits

Smaller $\longrightarrow$ Older visit

$0$ $\longrightarrow$ No visit

↳ SARSA

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha \Big[ R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \Big]$$

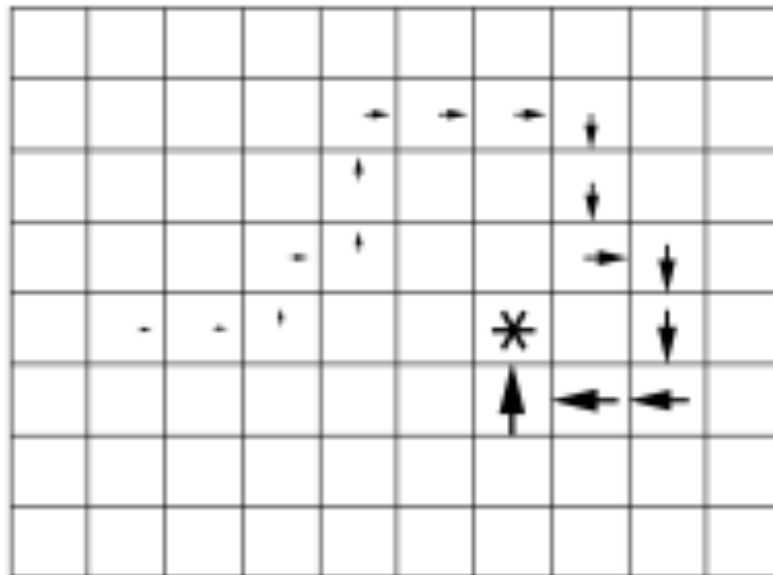observe $\rightarrow e(s_t, a_t) = e(s_t, a_t) + 1$

SARSA Error $\delta_t$

$$Q(s,a) = Q(s,a) + \alpha \, \delta_t \, e(s,a) \qquad \text{for all } s, a$$

next step $e(s_t, a_t) = \lambda \gamma e(s_t, a_t)$

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

■ Wall   ■ Bump   ■ Goal

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

■ Wall    ▨ Bump   ▨ Goal

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

■ Wall   ▨ Bump   ▨ Goal

$\uparrow$

## SARSA $-\lambda$

## episode 1

Random initial state

$(13)$   $R = \pi^{\varepsilon-g}(13) = \begin{cases} 0.25 \to U \ D \ R \ L \end{cases}$   $(13)$   $D = \pi^{\varepsilon-g}(13) = \begin{cases} 0.25 \ U \ D \ R \ L \end{cases}$

Reward $= -1$

$e(13, R) = e(13, R) + 1 = 0 + 1 = 1$

$\delta_t = R_{t+1} + \gamma \ Q(S_{t+1}, a_{t+1}) - Q(S_t, a_t)$

$\quad = -1 + 0.9 \times \underbrace{Q(13, D)}_{0} - \underbrace{Q(13, R)}_{0} = -1$

$Q(13, R) = \underline{Q(13, R)} + \underline{\alpha} \ \underline{\delta_t} \ \underline{e(13, R)} = -0.5$

$$e(13, R) = \underbrace{e(13, R)}_{1} \underbrace{\gamma}_{0.9} \underbrace{\lambda}_{0.95} = 0.855$$

⑬ —→ $D$ —————— ⑮ ———— $D = \pi^{\Xi-y}(13) = \begin{cases} 0.25 & \forall DRL \end{cases}$
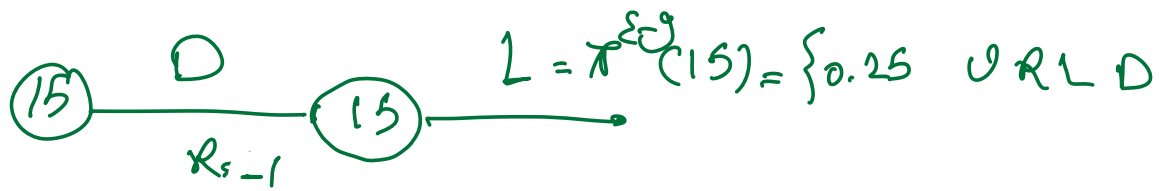
$R_{t+1} = -1$

$$e(13, D) = e(13, D) + 1 = 1$$

$$\delta_t = -1 + \gamma \underbrace{Q(15, L)}_{0} - \underbrace{Q(13, D)}_{0} = -1$$

$$Q(13, D) = \underbrace{Q(13, D)}_{0} + \underbrace{\alpha}_{0.5} \underbrace{\delta_t}_{-1} \underbrace{e(13, D)}_{1} = -0.5$$

$$Q(13, R) = \underbrace{Q(13, R)}_{-0.5} + \underbrace{\alpha}_{0.5} \underbrace{\delta_t}_{-1} \underbrace{e(13, R)}_{0.855} = -0.9275$$

$$e(13, D) = \overset{0.9}{\gamma} \overset{0.95}{\lambda} \overset{1}{e(13, D)} = 0.855$$

$$e(13, R) = \underbrace{\gamma}_{0.9} \underbrace{\lambda}_{0.75} \underbrace{e(13, R)}_{0.855} = 0.731$$

$$\overset{\textstyle 15}{\textcircled{15}} \xrightarrow{\quad D \quad} \overset{\textstyle }{\textcircled{15}} \xrightarrow{\quad} \qquad L = \pi^{\varepsilon \circ g}(15) = \begin{cases} 0.25 & U\,R\,L\,D \end{cases}$$

$R_{s-1}$

$$e(15, D) = e(15, D) + 1 = 1$$

$$\delta_t = -1 + \gamma\, Q(15, L) - Q(15, D) = -1$$

$$Q(15, D) = -0.5$$

$$Q(13, D) = Q(13, D) + \alpha\, \delta_t\, \underset{\underset{0.855}{\smile}}{e(13, D)} = -0.9275$$

$$Q(13, R) = Q(13, R) + \alpha\, \delta_t\, \underset{\underset{0.731}{\smile}}{e(13, R)} = -1.293$$

$$e(15, D) = \gamma \lambda\, e(15, D) = 0.855$$

$$e(13, D) = 0.731$$

$$e(13, R) = 0.625$$

---

$$\overset{\textstyle }{\textcircled{15}} \xrightarrow{\quad L \quad} \textcircled{Goal}$$

$R = 99$

$$e(15, L) = e(15, L) + 1 = 1$$

$$\delta_t = 99 + \gamma\, \underline{Q(Goal, a)} - \underline{Q(15, L)} = 99$$

$$Q(15, L) = Q(15, L) + \alpha \, \delta_t \, e(15, L) = 49.5$$

$$Q(15, D) = \underbrace{Q(15, D)}_{0} + \alpha \underbrace{\delta_t}_{0.5} \underbrace{}_{99} \underbrace{e(15, D)}_{0.855} = 41.82$$

$$Q(13, D) = Q(13, D) + \alpha \, \delta_t \, e(13, D) = 35.257$$

$$Q(13, R) = Q(13, R) + \alpha \, \delta_t \, e(13, R) = 29.84$$

## SARSA ($\lambda$) Algorithm

Initialize $Q(s, a)$ arbitrarily and $e(s, a) = 0$, for all $s, a$
Repeat (for each episode):
    Initialize $s, a$   $\leftarrow e(s, a) = 0$ for all $s, a$
    Repeat (for each step of episode):
        Take action $a$, observe $r, s'$
        Choose $a'$ from $s'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $\delta \leftarrow r + \gamma Q(s', a') - Q(s, a)$
        $e(s, a) \leftarrow e(s, a) + 1$
        For all $s, a$:
            $Q(s, a) \leftarrow Q(s, a) + \alpha \delta e(s, a)$
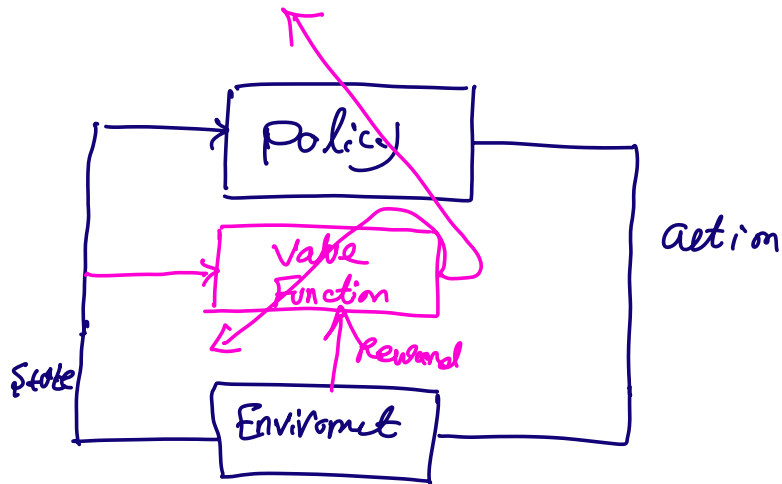            $e(s, a) \leftarrow \gamma \lambda e(s, a)$
        $s \leftarrow s'; a \leftarrow a'$
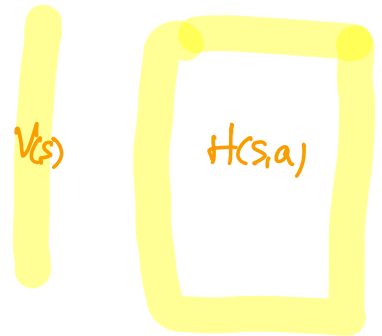    until $s$ is terminal

# Actor-critic Policy

$$\underset{a \in [\cdot \ 1]}{argmax} \ Q(s; a)$$

①-Learning & SARSA $\Rightarrow$ Value-based policy
$$Q \longrightarrow$$

$$V(s_t) = V(s_t) + \alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$
$$\delta_t$$
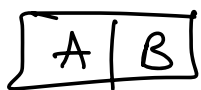


State

action

Reward

$V(s)$    $H(s,a)$

$$\delta_t = [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

$$V(s_t) = V(s_t) + \alpha \ \delta_t$$

$$H(s_t, a_t) = H(s_t, a_t) + \beta \ \delta_t \ (1 - \pi(a_t | s_t))$$
$\uparrow$
Preference

$$\pi(a_t | s_t) = \frac{e^{H(s_t, a_t)}}{\sum_{a \in A} e^{H(s_b, a)}}$$

$$\boxed{A \mid B}$$

$\alpha = 0.25 \quad , \quad \beta = 0.3$

$$\boxed{\begin{array}{l} H(s,a) = 0 \\ V(s) = 0 \end{array}}$$

$$M(a^1) = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$M(a^2) = \begin{bmatrix} \cdot & 1 \\ 1 & \cdot \end{bmatrix}$$

Reward $\begin{cases} B & +5 \\ a^2 & -1 \end{cases}$

## episode

Random State

$S_0 = A$

$a^2 = \pi(A) = \begin{cases} \pi(a^1 \mid A) = \dfrac{e^{H(A,a^1)}}{e^{H(A,a^1)} + e^{H(A,a^2)}} = \dfrac{e^0}{e^0 + e^0} = 0.5 \quad a^1 \\[4mm] \pi(a^2 \mid A) = \dfrac{e^{H(A,a^2)}}{e^{H(A,a^1)} + e^{H(A,a^2)}} = \dfrac{e^0}{e^0 + e^0} = 0.5 \quad a^2 \end{cases}$

$\xrightarrow[R = 4]{} B$

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t) = 4 + \gamma \cdot \underbrace{V(B)}_{0} - \underbrace{V(A)}_{0} = 4$$

$$\begin{cases} V(\underset{A}{S_t}) = V(\underset{A}{S_t}) + \underset{0.25}{\alpha} \underset{4}{\delta_t} = 1 \\[6mm] H(A,a^2) = \underbrace{H(A,a^2)}_{0} + \underset{0.3}{\beta} \times \underset{4}{\delta} \left( 1 - \underbrace{\pi(a^2 \mid A)}_{0.5} \right) = 0.6 \end{cases}$$

$B \xrightarrow[R=-1]{\pi =} A \quad \begin{cases} \pi(a^1 \mid B) = \dfrac{e^{H(B,a^1)}}{e^{H(B,a^1)} + e^{H(B,a^2)}} = \dfrac{e^0}{e^0 + e^0} = 0.5 \\[4mm] \pi(a^2 \mid B) = 1 - \pi(a^1 \mid B) = 0.5 \quad \leadsto a^2 \end{cases}$

$$\delta_t = -1 + \gamma V(A) - V(B) = -0.1$$
$$\underbrace{\phantom{V(A)}}_{1} \quad \underbrace{\phantom{V(B)}}_{0}$$

$$V(B) = \underbrace{V(B)}_{0} + \underbrace{\alpha}_{0.25} \underbrace{\delta_t}_{-0.1} = -0.025 \checkmark$$

$$H(B, a^2) = \underbrace{H(B, a^2)}_{0} + \underbrace{\beta}_{0.3} \underbrace{\delta_t}_{-0.1} (1 - \underbrace{\pi(a^2|B)}_{0.5}) = -0.015$$

$$A \xrightarrow[R=4]{\pi(a|A) = } B \quad \begin{cases} \pi(a'|A) = \dfrac{e^{H(A, a')}}{e^{H(A, a')} + e^{H(A, a^2)}} = \dfrac{e^0}{e^0 + e^{0.6}} = 0.3543 \\[2em] \pi(a^2|A) = 0.6457 \end{cases} \qquad \sim a^2$$

$$\delta = \underbrace{R}_{4} + \underbrace{\gamma}_{0.9} \underbrace{V(B)}_{-0.025} - \underbrace{V(A)}_{1} = 2.97$$

$$\begin{cases} V(A) = V(A) + \alpha \delta_t = 1 + 0.25 \times 2.87 = 1.74 \\[1em] H(A, a^2) = H(A, a^2) + \beta \delta (1 - \underbrace{\pi(a^2|A)}_{0.6457}) = 0.91 \end{cases}$$

---

$$B \xrightarrow[R=5]{} B \quad \begin{cases} \pi(a'|B) = 0.5037 \\ \pi(a^2|B) = 0.4963 \end{cases} \qquad \sim a'$$

$$\delta = \underbrace{R}_{5} + \underbrace{\gamma}_{0.9} \underbrace{V(B)}_{-0.025} - V(B) = 5.005$$

$$\begin{cases} V(B) = \underbrace{V(B)}_{-0.025} + \underbrace{\alpha}_{0.25} \underbrace{\delta_t}_{5.005} = 1.225 \\[1em] H(B, a') = 0.74 \end{cases}$$

$$H(A, a^1) = 0 \qquad H(B, a^1) = 0.74$$

$$H(A, a^2) = 0.96 \qquad H(B, a) = -0.015$$

$$\pi(a^1 | A) = \frac{e^{H(A, a^1)}}{e^{H(A, a^1)} + e^{H(A, a^1)}} = 0.28 \quad\Bigg| \quad \pi(a^1 | B) = 0.68$$

$$\pi(a^2 | A) = 0.72 \qquad\qquad\qquad \pi(a^2 | B) = 0.31$$

## Tabular Actor-Critic Algorithm

$V(s) = 0$, $H(s,a) = 0$, for all $s \in S$, $a \in A$.

Repeat for $N$ episodes

- Start from a random state $s_0 \in S$, $t = 0$

  While $t < T$ (episode Length).

  - Select action: $a_t \sim \pi(\cdot | s_t)$    $\pi(a|s) = \dfrac{e^{H(s,a)}}{\sum\limits_{a' \in A} e^{H(s,a')}}$

  - Take action $a_t$, move to state $s_{t+1}$ and observe $R_{t+1}$.

  - $\delta_t = R_{t+1} + \gamma\, V(s_{t+1}) - V(s_t)$

  - $V(s_t) = V(s_t) + \alpha\, \delta_t$

  - $H(s_t, a_t) = H(s_t, a_t) + \beta\, \delta_t\, (1 - \pi(a_t | s_t))$

  - $t = t+1$