# COMPSCI 389
# Introduction to Machine Learning

**Value functions, Temporal Difference Learning, and Actor-Critics**

Prof. Philip S. Thomas (pthomas@cs.umass.edu)

# REINFORCE (Review)

---

**Algorithm 17.2:** REINFORCE

---

1 **for** *each episode* **do**

2     // Run one episode (play one game).

3     **for** *each time t in the episode* **do**

4        Agent observes state $S_t$;

5        Agent selects action $A_t$ according to the current policy, $\pi_\theta$;

6        Environment responds by transitioning from state $S_t$ to state $S_{t+1}$ and emitting reward $R_t$;

7     **end**

8     // Learn from the outcome of the episode.

9     **for** *each time t in the episode* **do**

10        $\forall i, \; \theta_i \leftarrow \theta_i + \alpha\gamma^t \left(\sum_{k=0}^{\infty} \gamma^k R_{t+k}\right) \frac{\partial \ln(\pi_\theta(S_t, A_t))}{\partial \theta_i}$;

11     **end**

12 **end**

---

2

# Two-Phases

- REINFORCE has two phases:
  - Run an episode, selecting actions and observing the outcome.
  - Update the policy.
- Waiting until the end of an episode to update the policy seems inefficient.

**Algorithm 17.2: REINFORCE**

1 **for** *each episode* **do**
2     // Run one episode (play one game).
3     **for** *each time t in the episode* **do**
4         Agent observes state $S_t$;
5         Agent selects action $A_t$ according to the current policy, $\pi_\theta$;
6         Environment responds by transitioning from state $S_t$ to state $S_{t+1}$ and emitting reward $R_t$;
7     **end**
8     // Learn from the outcome of the episode.
9     **for** *each time t in the episode* **do**
10         $\forall i, \theta_i \leftarrow \theta_i + \alpha \gamma^t \left( \sum_{k=0}^{\infty} \gamma^k R_{t+k} \right) \frac{\partial \ln(\pi_\theta(S_t, A_t))}{\partial \theta_i}$;
11     **end**
12 **end**

# Idea: Consider expected future rewards

- Imagine that you play the lottery and learn that you have won.

- In the future, you will obtain the prize money, which may result in actual rewards (in the form of increased comfort and pleasure).

- However, you will likely be happy and celebrate, perhaps learning that you should play the lottery more, *before* you collect any prize money or see the actual impact that it has on your life!

- This isn't due to any rewards – you have received no actual rewards yet.
    - This learning (change in behavior) is due to your expectations of future rewards.

# Idea: Consider expected future rewards

- When you realize that you have won the lottery, you recognize that the expected outcome was better than you were previously expecting.
  - This is sufficient for you to learn!

# Idea: Consider expected future rewards

- Consider another example: a monkey

Firing of neurons that indicate "take recent actions more often."

(Details later!)

Time

R

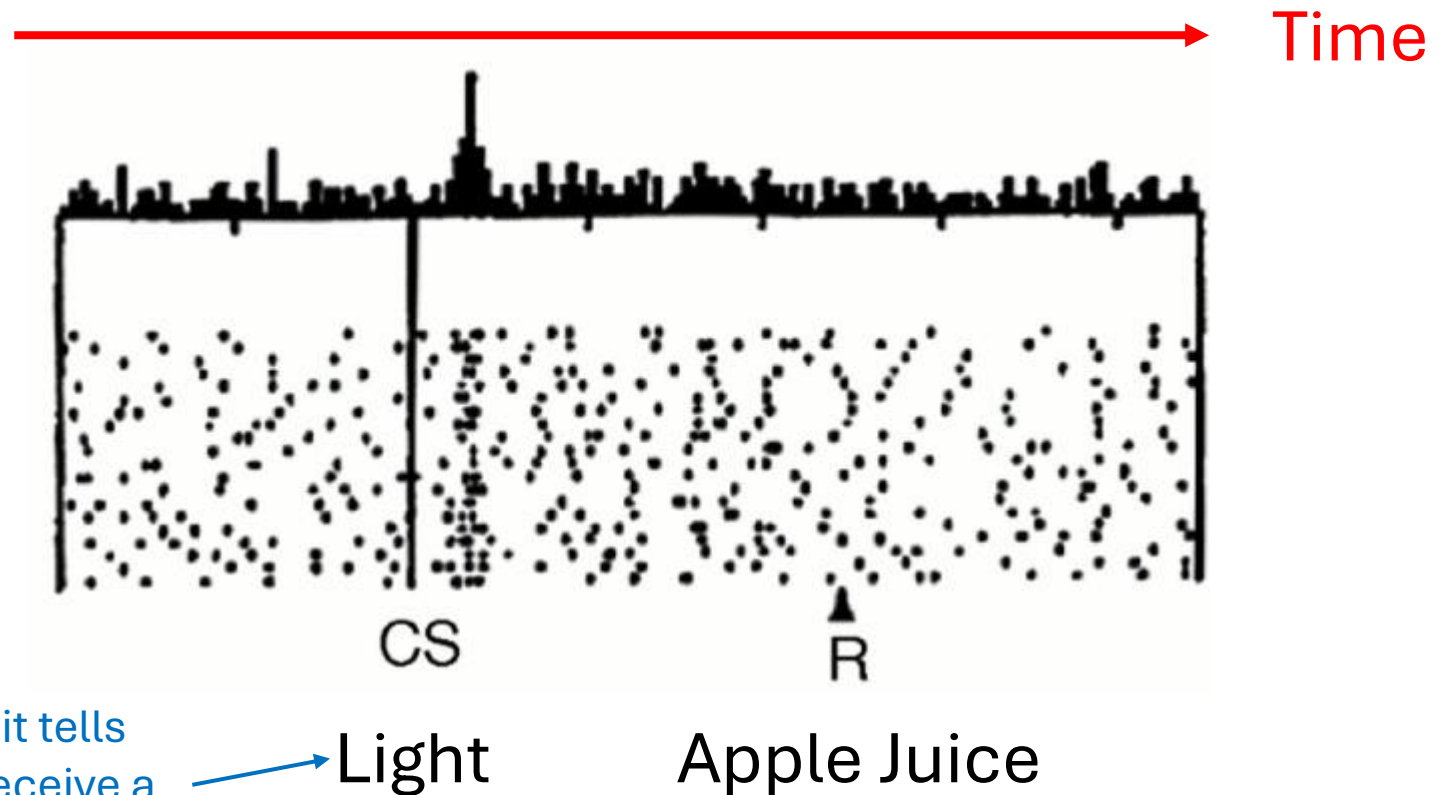Apple Juice

# Idea: Consider expected future rewards

• Consider another example: a monkey



Time

Firing of neurons that indicate "take recent actions more often."

(Details later!)

CS

R

Light          Apple Juice
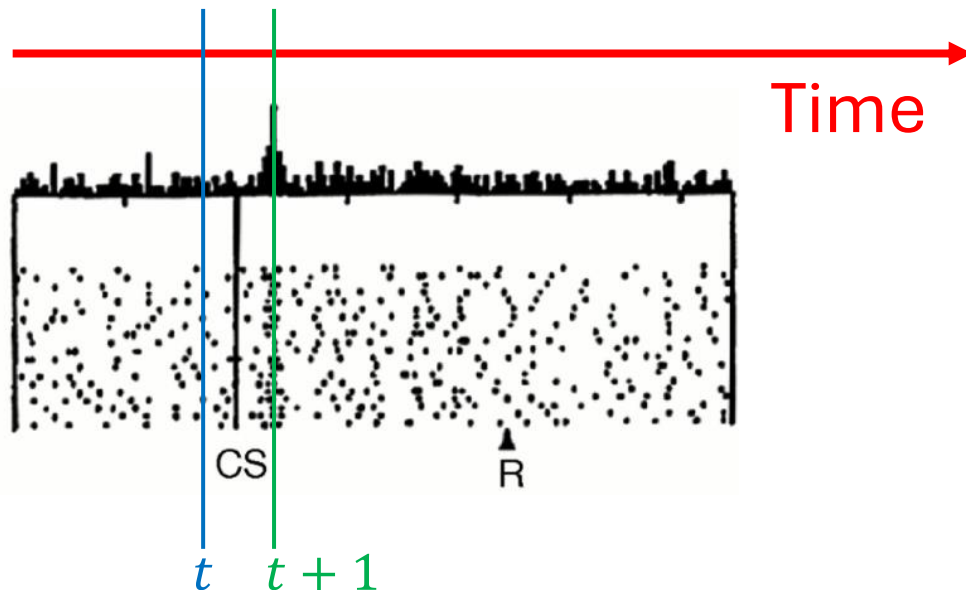
Seeing a light flash isn't a reward, but it tells the agent (monkey) that it is likely to receive a reward (apple juice) in the future.

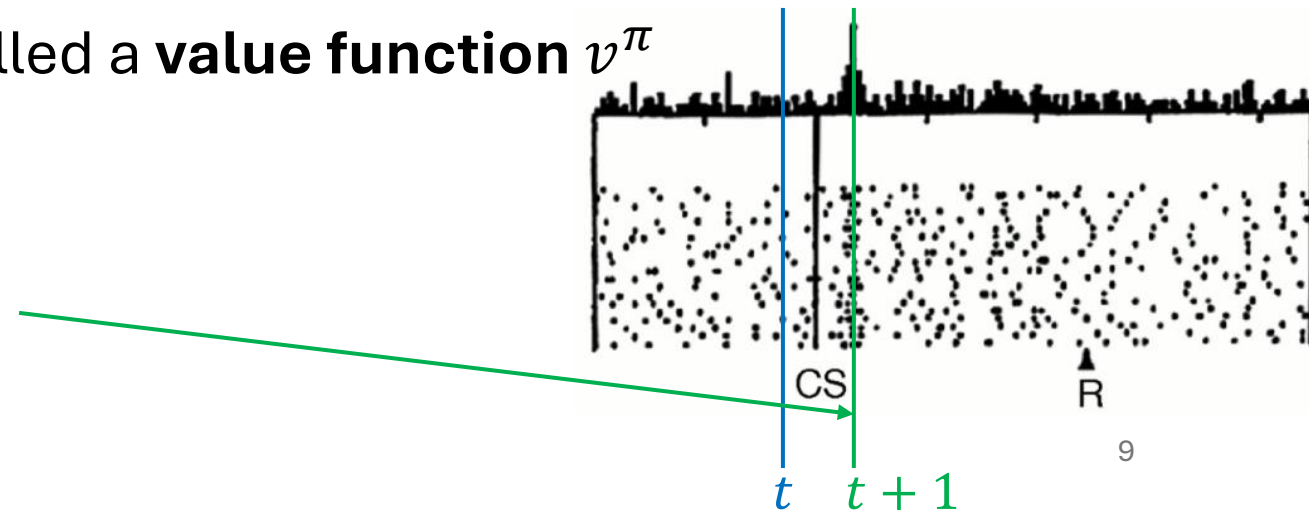# Idea: Consider expected future rewards

- **Before**: Make actions more likely when they result in an observed desirable outcome

- **New idea**: Make actions more likely when they cause the agent to believe that the outcome will be more desirable than expected.

# Idea: Consider expected future rewards

- **Before**: Make actions more likely when they result in an observed desirable outcome

- **New idea**: Make actions more likely when they cause the agent to believe that the outcome will be more desirable than expected.

- To do this, the agent must have a notion of how much reward it expects to get in the future.
    - This is captured by a function called a **value function** $v^\pi$

To know that recent actions should be reinforced, we must know that seeing the light means that more reward is expected in the future.

CS

R

$t$   $t+1$

9

# Value Function $v^{\pi}$

- A function that depends on the policy $\pi$
  - **Input**: Any state $s$
  - **Output**: The expected discounted sum of rewards that the agent will receive in the future if it were in state $s$:

$$v^{\pi}(s) = \mathbf{E}\left[\sum_{k=0}^{\infty} \gamma^t R_{t+k} \middle| S_t = s; \pi\right].$$

- The right-hand side has $t$, but $t$ doesn't show up on the left side!
  - Due to the **Markov property**, the value of $v^{\pi}(s)$ is the same for all $t$ on the right-hand side.

# Value Function $v^\pi$

- Consider a gridworld
  - Deterministic transitions
  - The agent always starts in the top-left.
  - The goal (a terminal state) is the bottom-right.
  - The reward is 0 at each time step, except when the agent **enters** the goal state, at which time it receives a reward of 100.
  - Let $\gamma = 0.5$
- Let $\pi$ be an optimal policy. The value function is then$\rightarrow$

|  |  |  |  |  |
|---|---|---|---|---|
|  |  |  |  |  |
|  |  |  |  |  |
|  |  |  |  |  |
|  |  |  |  |  |
|  |  |  |  | ? |

# Value Function $v^\pi$

- Consider a gridworld
  - Deterministic transitions
  - The agent always starts in the top-left.
  - The goal (a terminal state) is the bottom-right.
  - The reward is 0 at each time step, except when the agent **enters** the goal state, at which time it receives a reward of 100.
  - Let $\gamma = 0.5$
- Let $\pi$ be an optimal policy. The value function is then→

| | | | | |
|---|---|---|---|---|
| | | | | |
| | | | | |
| | | | | |
| | | | | ? |
| | | | | 0 |

# Value Function $v^{\pi}$

- Consider a gridworld
  - Deterministic transitions
  - The agent always starts in the top-left.
  - The goal (a terminal state) is the bottom-right.
  - The reward is 0 at each time step, except when the agent **enters** the goal state, at which time it receives a reward of 100.
  - Let $\gamma = 0.5$
- Let $\pi$ be an optimal policy. The value function is then→

| | | | | |
|---|---|---|---|---|
| | | | | |
| | | | | |
| | | | | ? |
| | | | | 100 |
| | | | | 0 |

# Value Function $v^{\pi}$

- Consider a gridworld
  - Deterministic transitions
  - The agent always starts in the top-left.
  - The goal (a terminal state) is the bottom-right.
  - The reward is 0 at each time step, except when the agent **enters** the goal state, at which time it receives a reward of 100.
  - Let $\gamma = 0.5$
- Let $\pi$ be an optimal policy. The value function is then→

| | | | | |
|---|---|---|---|---|
| | | | | |
| | | | | |
| | | | | 50 |
| | | | | 100 |
| | | | | 0 |

# Value Function $v^\pi$

- Consider a gridworld
  - Deterministic transitions
  - The agent always starts in the top-left.
  - The goal (a terminal state) is the bottom-right.
  - The reward is 0 at each time step, except when the agent **enters** the goal state, at which time it receives a reward of 100.
  - Let $\gamma = 0.5$
- Let $\pi$ be an optimal policy. The value function is then$\rightarrow$

| $\approx 0.78$ | 1.5625 | 3.125 | 6.25 | 12.5 |
|---|---|---|---|---|
| 1.5625 | 3.125 | 6.25 | 12.5 | 25 |
| 3.125 | 6.25 | 12.5 | 25 | 50 |
| 6.25 | 12.5 | 25 | 50 | 100 |
| 12.5 | 25 | 50 | 100 | 0 |

# Value Function $v^\pi$

So, different policies result in different value functions!

- Consider a gridworld
  - Deterministic transitions
  - The agent always starts in the top-left.
  - The goal (a terminal state) is the bottom-right.
  - The reward is 0 at each time step, except when the **enters** the goal state, at which time it receives a reward of 100.
  - Let $\gamma = 0.5$

- Let $\pi$ be the policy that always selects the down action. The value function is then$\rightarrow$

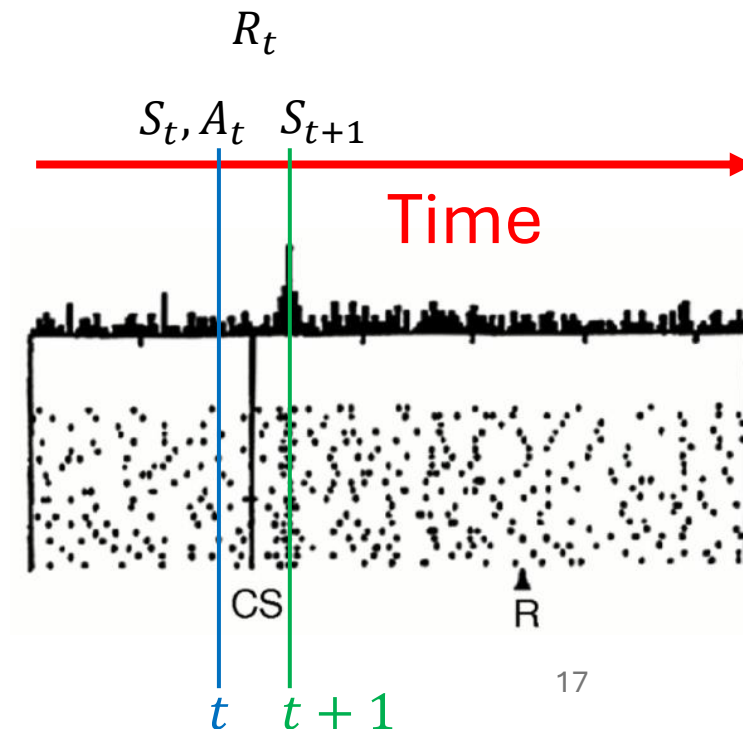| 0 | 0 | 0 | 0 | 12.5 |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 25 |
| 0 | 0 | 0 | 0 | 50 |
| 0 | 0 | 0 | 0 | 100 |
| 0 | 0 | 0 | 0 | 0 |

# Temporal Difference Error

- Later we will discuss how the agent can learn (estimate) $v^\pi$ from its experiences.

- First, let's explore how the agent can *use* $v^\pi$ (or an estimate of $v^\pi$).

- Consider how an agent can update its policy when:
  - The agent observes $S_t$
  - The agent selects $A_t$
  - The environment transitions to $S_{t+1}$ and emits reward $R_t$

$v^\pi(S_t)$: Expected total future reward

$R_t$: Reward received

$v^\pi(S_{t+1})$: Expected total future reward at next time

$$\underline{\hspace{2cm}} \approx R_t + \underline{\hspace{2cm}}$$

$R_t$

$S_t, A_t \quad S_{t+1}$

Time

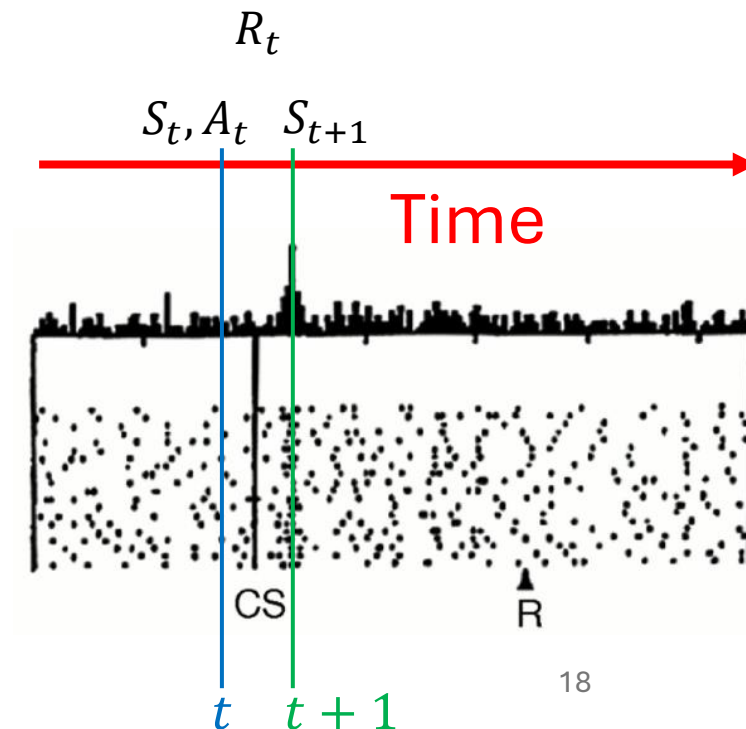CS        R

$t \quad t+1$

# Temporal Difference Error

- Later we will discuss how the agent can learn (estimate) $v^\pi$ from its experiences.

- First, let's explore how the agent can *use* $v^\pi$ (or an estimate of $v^\pi$).

- Consider how an agent can update its policy when:
  - The agent observes $S_t$
  - The agent selects $A_t$
  - The environment transitions to $S_{t+1}$ and emits reward $R_t$

$v^\pi(S_t)$: Expected total future reward

$R_t$: Reward received

$v^\pi(S_{t+1})$: Expected total future reward at next time

$$v^\pi(S_t) \approx R_t + v^\pi(S_{t+1})$$

# Temporal Difference Error

- So, we expect that $v^\pi(S_t) \approx R_t + v^\pi(S_{t+1})$
  - Note: With reward discounting, $v^\pi(S_t) \approx R_t + \gamma v^\pi(S_{t+1})$

- **Question**: What does it mean if $R_t + \gamma v^\pi(S_{t+1}) > v^\pi(S_t)$?



Time

$t \quad t+1$

$t \quad t+1$

This can happen when the reward is bigger than expected

This can happen when both

This can happen when the next state is better than expected

# Temporal Difference Error

- So, we expect that $v^\pi(S_t) \approx R_t + v^\pi(S_{t+1})$
  - Note: With reward discounting, $v^\pi(S_t) \approx R_t + \gamma v^\pi(S_{t+1})$
- **Question**: What does it mean if $R_t + \gamma v^\pi(S_{t+1}) > v^\pi(S_t)$?
- **Answer**: The outcome of $A_t$ was better than expected.
  - Make $A_t$ *more* likely in state $S_t$
- **Question**: What does it mean if $v^\pi(S_t) > R_t + \gamma v^\pi(S_{t+1})$?
- **Answer**: The outcome of $A_t$ was *worse* than expected.
  - Make $A_t$ *less* likely in state $S_t$

# Temporal Difference Error

- Let the **temporal difference error** or **TD error** $\delta_t$ be:
$$\delta_t = R_t + \gamma v^\pi(S_{t+1}) - v^\pi(S_t)$$
- When the TD error is positive, $R_t + \gamma v^\pi(S_{t+1}) > v^\pi(S_t)$
  - The agent should select $A_t$ more often in $S_t$.
  - The magnitude of the TD error indicates how much better the outcome of $A_t$ was, and can be used to scale how much more likely $A_t$ is made.
- When the TD error is negative, $v^\pi(S_t) > R_t + \gamma v^\pi(S_{t+1})$
  - The agent should select $A_t$ less often in $S_t$.
  - The magnitude [...]
- Idea: Replace the weight, $\displaystyle\sum_{t'=t}^{\infty} \gamma^{t'} R_{t'} = \gamma^t \sum_{k=0}^{\infty} \gamma^k R_{t+k}$ with $\delta_t$.
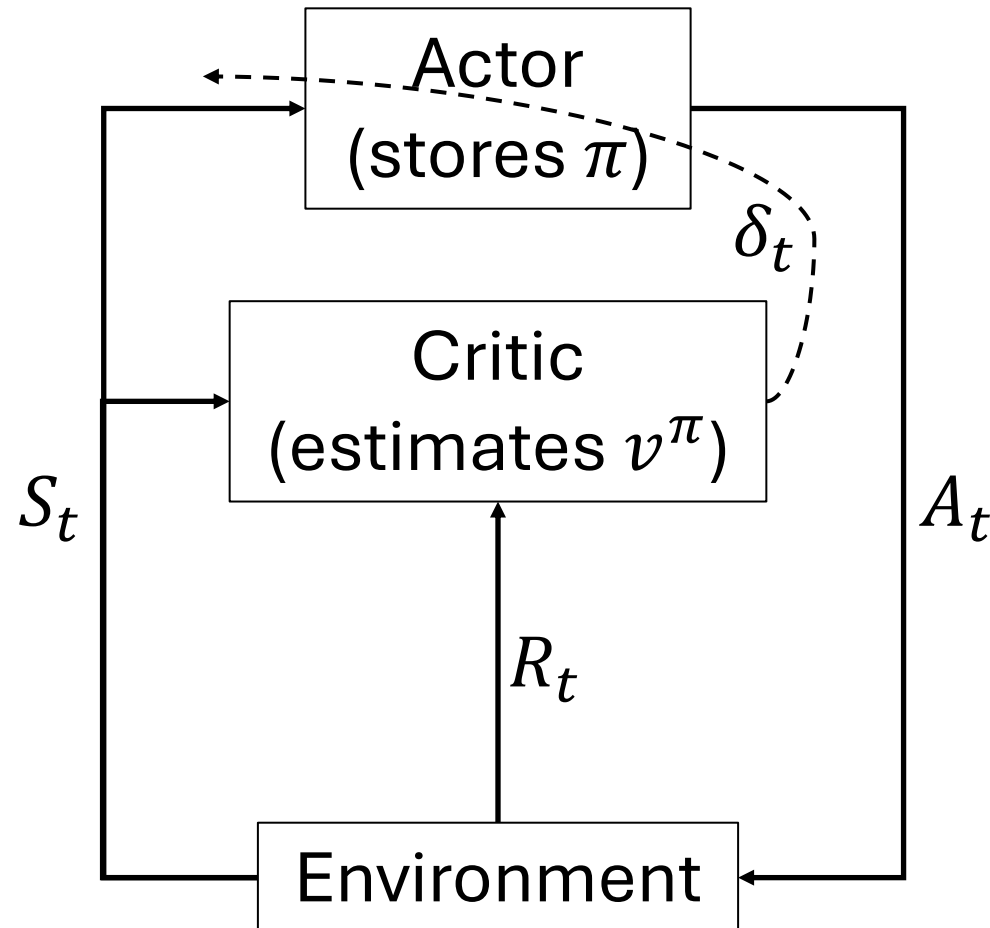
**Algorithm 17.3:** A simple RL algorithm inspired by MENACE, Version 5.0

---

1 **for** *each episode* **do**
2      // Run one episode (play one game).
3      **for** *each time t in the episode* **do**
4          // Execute one time step of agent-environment
                interaction
5          Agent observes state $S_t$;
6          Agent selects action $A_t$ according to the current policy, $\pi_\theta$;
7          Environment responds by transitioning from state $S_t$ to state $S_{t+1}$ and emitting reward $R_t$;
8          // Learn from the outcome of this one time step
9          $\delta_t \leftarrow R_t + \gamma v^\pi(S_{t+1}) - v^\pi(S_t)$;
10          $\forall i, \ \theta_i \leftarrow \theta_i + \alpha \delta_t \frac{\partial \pi_\theta(S_t, A_t)}{\partial \theta_i}$;
11      **end**
12 **end**

**Note:** We have replaced the discounted sum of rewards after $A_t$ with the TD error. This allows us to perform policy updates *before* the episode ends!

# Actor-Critic

$$\delta_t \leftarrow R_t + \gamma v^\pi(S_{t+1}) - v^\pi(S_t);$$
$$\forall i, \ \theta_i \leftarrow \theta_i + \alpha \delta_t \frac{\partial \pi_\theta(S_t, A_t)}{\partial \theta_i};$$

# Actor-Critics

- Actor-critic algorithms are a class of algorithms, not one specific algorithm.
  - They have an **actor** (which stores the current policy) and a **critic** (which stores an estimate of a value function).
- Often actor-critic algorithms are **policy gradient** algorithms like REINFORCE: they change the policy following **estimates** of the gradient of

$$J(\theta) = \mathbf{E}\left[\sum_{t=0}^{\infty} \gamma^t R_t\right]$$

  - Usually these estimates are **biased** (not only do they have variance, but even with infinite data they point in the "wrong" direction!)
  - Still, these biased estimates can be quite effective!
    - State of the art algorithms like PPO and SAC are actor-critics / policy gradient methods.

# Approximating the value function

- Let $v_w$ be a parametric function approximator (like a model in supervised learning) with weights $w$.

- Our aim is to make $v_w(s) \approx v^\pi(s) = \mathbf{E}[\sum_{k=0}^\infty \gamma^k R_t \mid S_t = s]$.

- Recall: $\delta_t = R_t + \gamma v^\pi(S_{t+1}) - v^\pi(S_t)$

- We don't know $v^\pi$, so let's re-define the TD error to use the approximation:
$$\delta_t = R_t + \gamma v_w(S_{t+1}) - v_w(S_t).$$

# Learning the value function

- Recall: $\delta_t = R_t + \gamma v_w(S_{t+1}) - v_w(S_t)$.
- **Question**: If $\delta_t > 0$, what does that say about $v_w(S_t)$?
- **Answer**: It should be increased!
- **Question**: If $\delta_t < 0$, what does that say about $v_w(S_t)$?
- **Answer**: It should be decreased!
- **Question**: How do we change $w$ to increase $v_w(S_t)$?
- **Answer**: $w \leftarrow w + \alpha \dfrac{\partial v_w(S_t)}{\partial w}$
- **Update**: $w \leftarrow w + \alpha \delta_t \dfrac{\partial v_w(S_t)}{\partial w}$

## Algorithm 18.2: An Actor-Critic Algorithm

**1 for** *each episode* **do**
**2**    // Run one episode (play one game).
**3**    **for** *each time t in the episode* **do**
**4**      // Execute one time step of agent-environment
       interaction
**5**      Agent observes state $S_t$;
**6**      Agent selects action $A_t$ according to the current policy, $\pi_\theta$;
**7**      Environment responds by transitioning from state $S_t$ to state $S_{t+1}$ and emitting reward $R_t$;
**8**      // Learn from the outcome of this one time step
**9**      $\delta_t \leftarrow R_t + \gamma v_w(S_{t+1}) - v_w(S_t)$;
**10**     $\forall i, \theta_i \leftarrow \theta_i + \alpha \delta_t \frac{\partial \pi_\theta(S_t, A_t)}{\partial \theta_i}$ // Actor update
**11**     $\forall j, w_j \leftarrow w_j + \beta \delta_t \frac{\partial v_w(S_t)}{\partial w_j}$ // Critic update
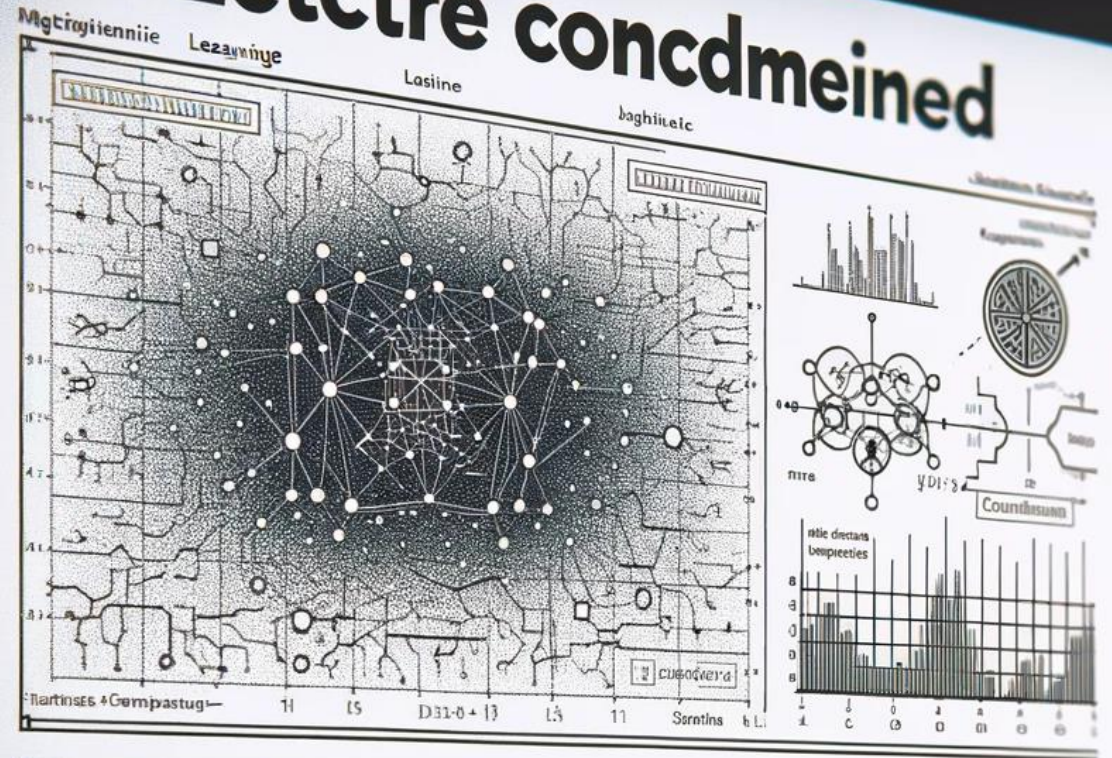**12**    **end**
**13 end**

Compute the TD error, $\delta_t$.

Update the actor's policy. Optionally, use $\partial \ln(\pi_\theta(S_t, A_t))$.

Update the critic's value function estimate. Optionally, use a step size $\beta$ that is different from that of the actor ($\alpha$).

27

End