

---

# Model Free Reinforcement Learning

## $\epsilon$ -Greedy Exploration

## Q-Learning

## Function Approximation

---

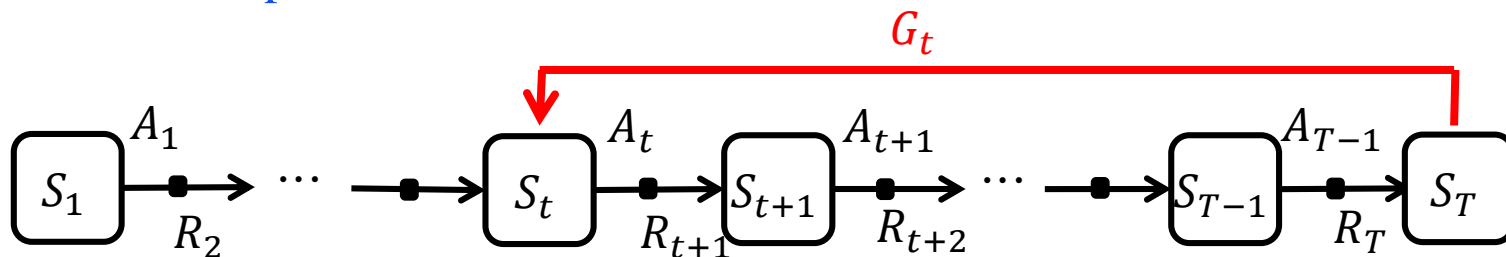
# Model Free Reinforcement Learning

# Model Free Reinforcement Learning

- No model
- Learn value function (and/or policy) from experience
- Common Model Free RL
  - Monte-Carlo (MC) Reinforcement Learning
  - Temporal Difference (TD) Reinforcement Learning
  - TD( $\lambda$ )

# Monte-Carlo Reinforcement Learning

- MC methods learn **directly from episodes of experience**
- MC is **model-free**:
  - no knowledge of MDP transitions / rewards
- MC learns from **complete episodes**:
  - no bootstrapping
- MC uses the simplest possible idea:
  - **value = mean return**
- Caveat: can only apply MC to episodic MDPs
  - **All episodes must terminate**



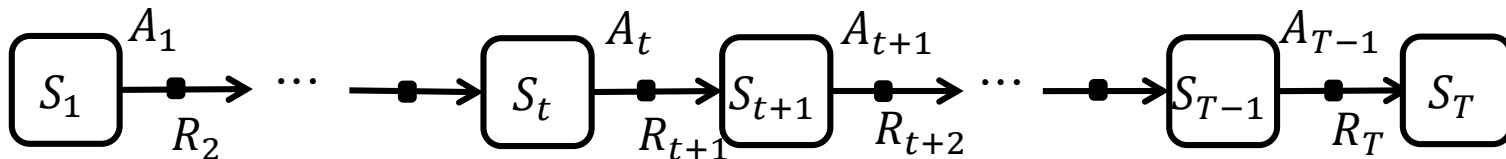
# Monte-Carlo Policy Evaluation

- Goal: learn  $v_\pi$  from episodes of experience under policy  $\pi$   
 $S_1, A_1, R_2, \dots, S_T \sim \pi$
- Recall that the return is the total discounted reward:  
$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$
- Recall that the value function is the expected return:  
$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$$
- Monte-Carlo policy evaluation uses empirical mean return instead of expected return



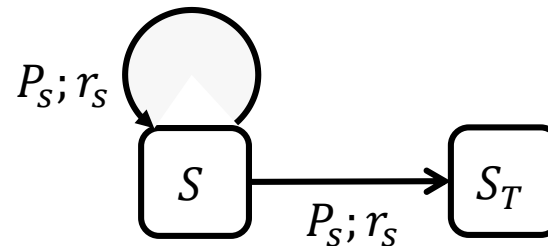
# Monte-Carlo Policy Evaluation (cont.)

- To evaluate  $v_\pi(s)$  at state  $s$ 
  - Increment counter  $N(s) \leftarrow N(s) + 1$
  - Increment total return  $S(s) \leftarrow S(s) + G_t$
  - Value is estimated by mean return  $V(s) \leftarrow S(s)/N(s)$
- By law of large numbers,  $V(s) \rightarrow v_\pi(s)$  as  $N(s) \rightarrow \infty$



# First Visit vs. Every Visit

- To evaluate  $v_\pi(s)$  at state  $s$ 
  - Increment counter  $N(s) \leftarrow N(s) + 1$
  - Increment total return  $S(s) \leftarrow S(s) + G_t$
  - Value is estimated by mean return  $V(s) \leftarrow S(s)/N(s)$
- By law of large numbers,  $V(s) \rightarrow v_\pi(s)$  as  $N(s) \rightarrow \infty$
- What if the same state  $s$  is visited in an episode?
  - Do the above for **every visit or first visit**?
    - ▶ What happen for the case in the figure?
    - ▶ Both converge quadratically, so this issue is ignored in this course.



# Incremental Mean

The mean  $\mu_1, \mu_2, \dots$  of a sequence  $x_1, x_2, \dots$  can be computed incrementally,

$$\begin{aligned}\mu_k &= \frac{1}{k} \sum_{j=1}^k x_j \\ &= \frac{1}{k} (x_k + \sum_{j=1}^{k-1} x_j) \\ &= \frac{1}{k} (x_k + (k-1)\mu_{k-1}) \\ &= \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1})\end{aligned}$$



# Incremental Monte-Carlo Updates

- Update  $V(s)$  incrementally after episode  $S_1, A_1, R_2, \dots, S_T$
- For each state  $S_t$  with return  $G_t$

$$N(S_t) \leftarrow N(S_t) + 1$$

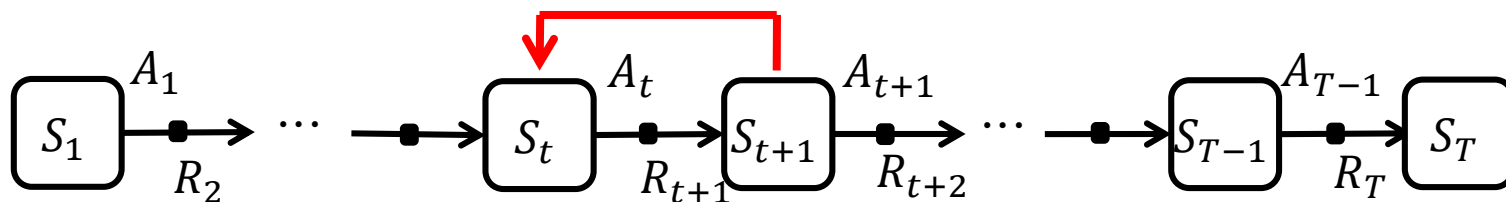
$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} (G_t - V(S_t))$$

- In **non-stationary problems**, it can be useful to track a running mean, i.e. **forget old episodes**.

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$

# Temporal-Difference Learning

- TD methods **learn directly from episodes of experience**
- TD is **model-free**:
  - no knowledge of MDP transitions / rewards
- TD learns from **incomplete episodes**,
  - by bootstrapping
- TD updates a **guess** towards a **guess**



# MC vs. TD

- Goal: learn  $v_\pi$  online from experience under policy  $\pi$
- Incremental every-visit Monte-Carlo
  - Update value  $V(S_t)$  toward actual return  $G_t$ 
$$V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$$
- Simplest temporal-difference learning algorithm: TD(0)
  - Update value  $V(S_t)$  toward estimated return  $R_{t+1} + \gamma V(S_{t+1})$ 
$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$
  - $R_{t+1} + \gamma V(S_{t+1})$  is called the **TD target**
  - $\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$  is called the **TD error**

## TD vs. MC (I)

- TD can learn **before** knowing the final outcome
  - TD can learn online after every step
  - MC must wait until end of episode before return is known
- TD can learn **without** the final outcome
  - TD can learn from incomplete sequences
  - MC can only learn from complete sequences
  - TD works in continuing (non-terminating) environments
  - MC only works for episodic (terminating) environments

## Bias/Variance Trade-Off

- TD target  $R_{t+1} + \gamma V(S_{t+1})$  is **biased estimate** of  $v_\pi(S_t)$ 
  - Return  $G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$  is **unbiased estimate** of  $v_\pi(S_t)$
  - True TD target  $R_{t+1} + \gamma v_\pi(S_{t+1})$  is **unbiased estimate** of  $v_\pi(S_t)$
- TD target is **much lower variance than the return**:
  - Return depends on **many** random actions, transitions, rewards
  - TD target depends on **only one** random action, transition, reward

## MC vs. TD (II)

- MC has high variance, zero bias
  - Good convergence properties (even with function approximation)
  - Not very sensitive to initial value
  - Very simple to understand and use
- TD has low variance, some bias
  - Usually more efficient than MC
  - TD(0) converges to  $v_{\pi}(s)$  (but not always with function approximation)
  - More sensitive to initial value

## Batch MC and TD

- MC and TD converge:  $V(s) \rightarrow v_\pi(s)$  as experience  $\rightarrow \infty$
- But what about batch solution for finite experience?

$$s_1^1, a_1^1, r_2^1, \dots, s_{T_1}^1$$

$$\vdots$$

$$s_1^k, a_1^k, r_2^k, \dots, s_{T_k}^k$$

- e.g. Repeatedly sample episode  $k \in [1, K]$
- Apply MC or TD(0) to episode  $k$

# AB Example

Two states A, B; no discounting; 8 episodes of experience

A, 0, B, 0

B, 1

B, 1

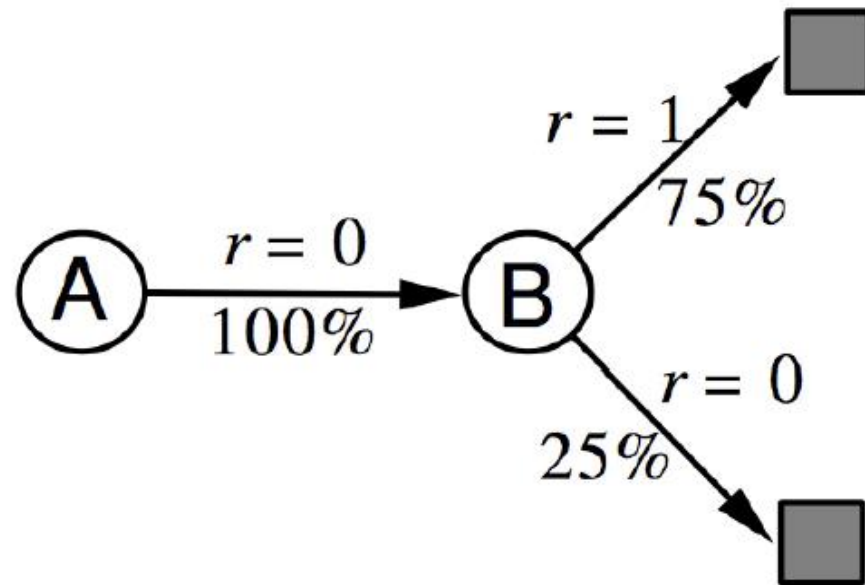
B, 1

B, 1

B, 1

B, 1

B, 0



What is  $V(A)$ ,  $V(B)$ ?

Both MC and TD will obtain different values!!



# Certainty Equivalence

- MC converges to solution with minimum mean-squared error
  - Best fit to the observed returns

$$\sum_{k=1}^K \sum_{t=1}^{T_k} (G_t^k - V(s_t^k))^2$$

- In the AB example,  $V(A) = 0$ ,  $V(B) = 0.75$

- TD(0) converges to solution of max likelihood Markov model
  - Solution to the MDP  $\langle \mathcal{S}, \mathcal{A}, \hat{\mathcal{P}}, \hat{\mathcal{R}}, \gamma \rangle$  that best fits the data

$$\hat{p}_s^a = \frac{1}{N(s, a)} \sum_{k=1}^K \sum_{t=1}^{T_k} 1(s_t^k, a_t^k, s_{t+1}^k = s, a, s')$$

$$\hat{\mathcal{R}}_s^a = \frac{1}{N(s, a)} \sum_{k=1}^K \sum_{t=1}^{T_k} 1(s_t^k, a_t^k = s, a) r_t^k$$

- In the AB example,  $V(A) = 0.75$ ,  $V(B) = 0.75$

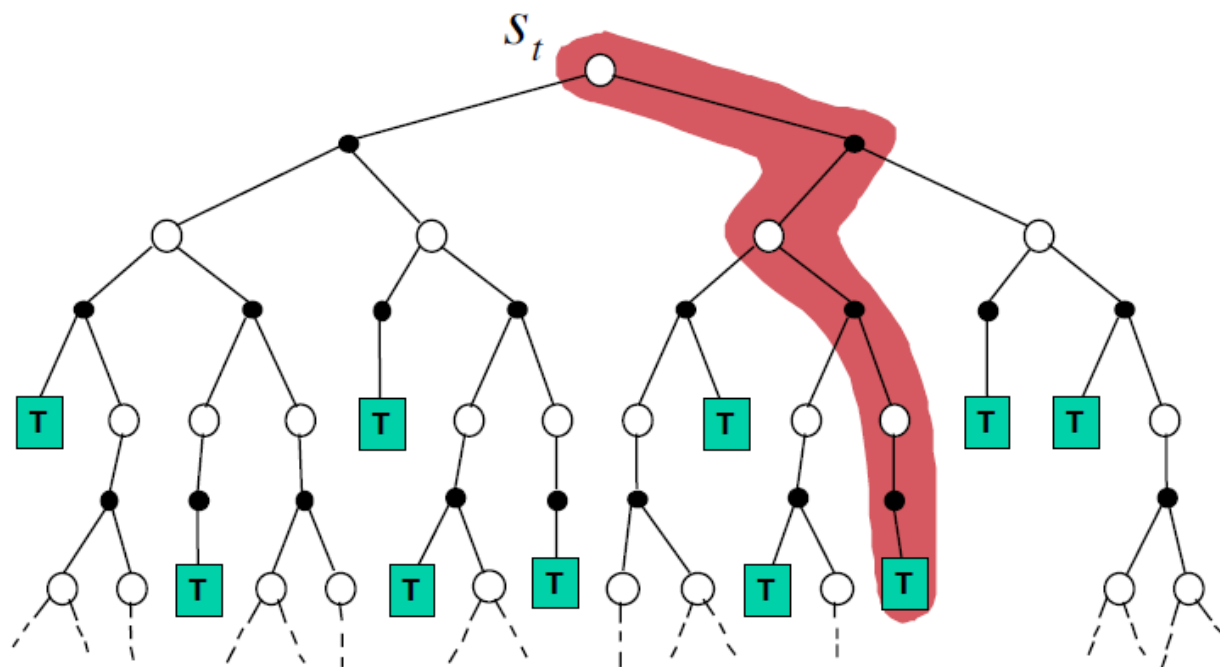


## MC vs. TD (III)

- **TD exploits Markov property**
  - Usually **more efficient in Markov environments**
    - ▶ So, TD works well for MDP problems like 2048.
- **MC does not exploit Markov property**
  - Usually **more effective in non-Markov environments**
    - ▶ MC works fine for non-MDP too.

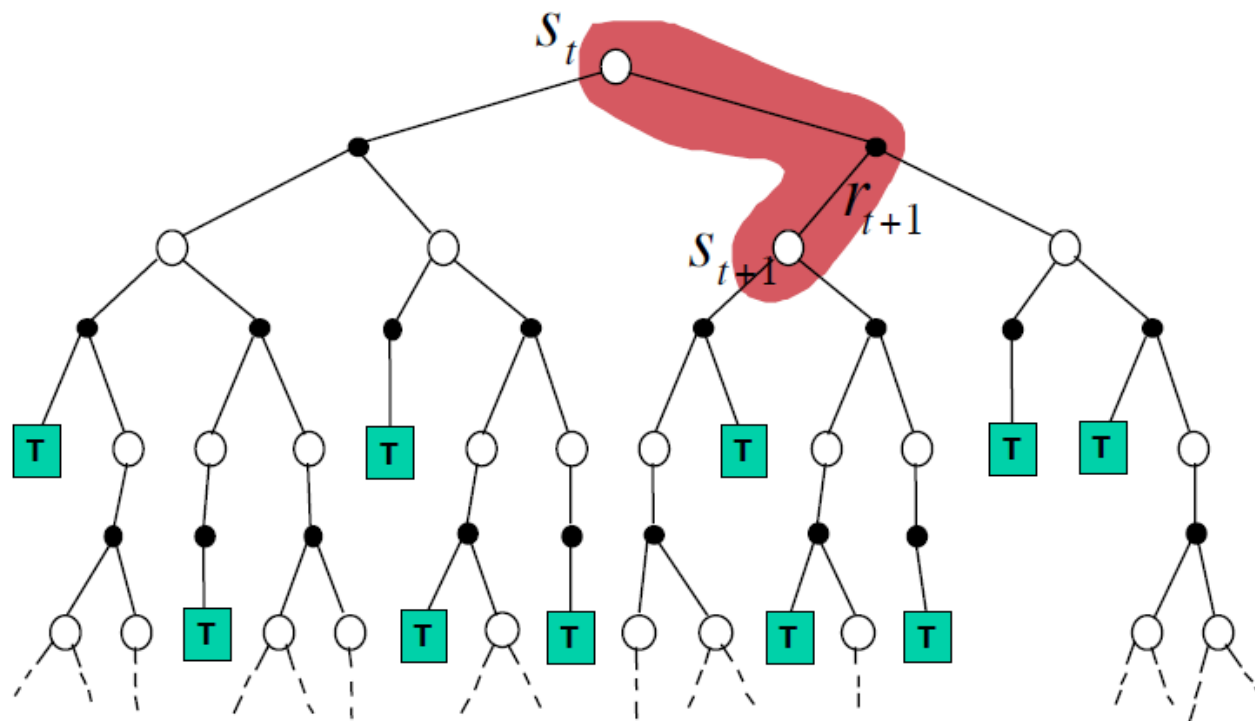
# Monte-Carlo Backup

$$V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$$



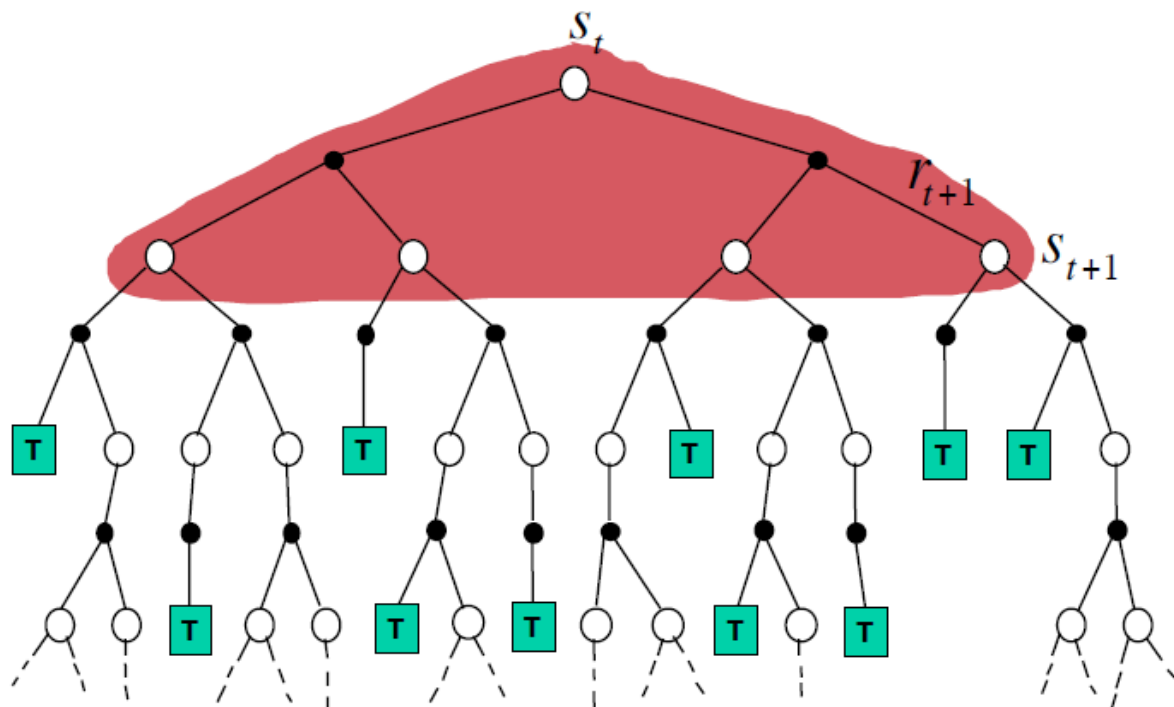
# Temporal-Difference Backup

$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



# Dynamic Programming Backup

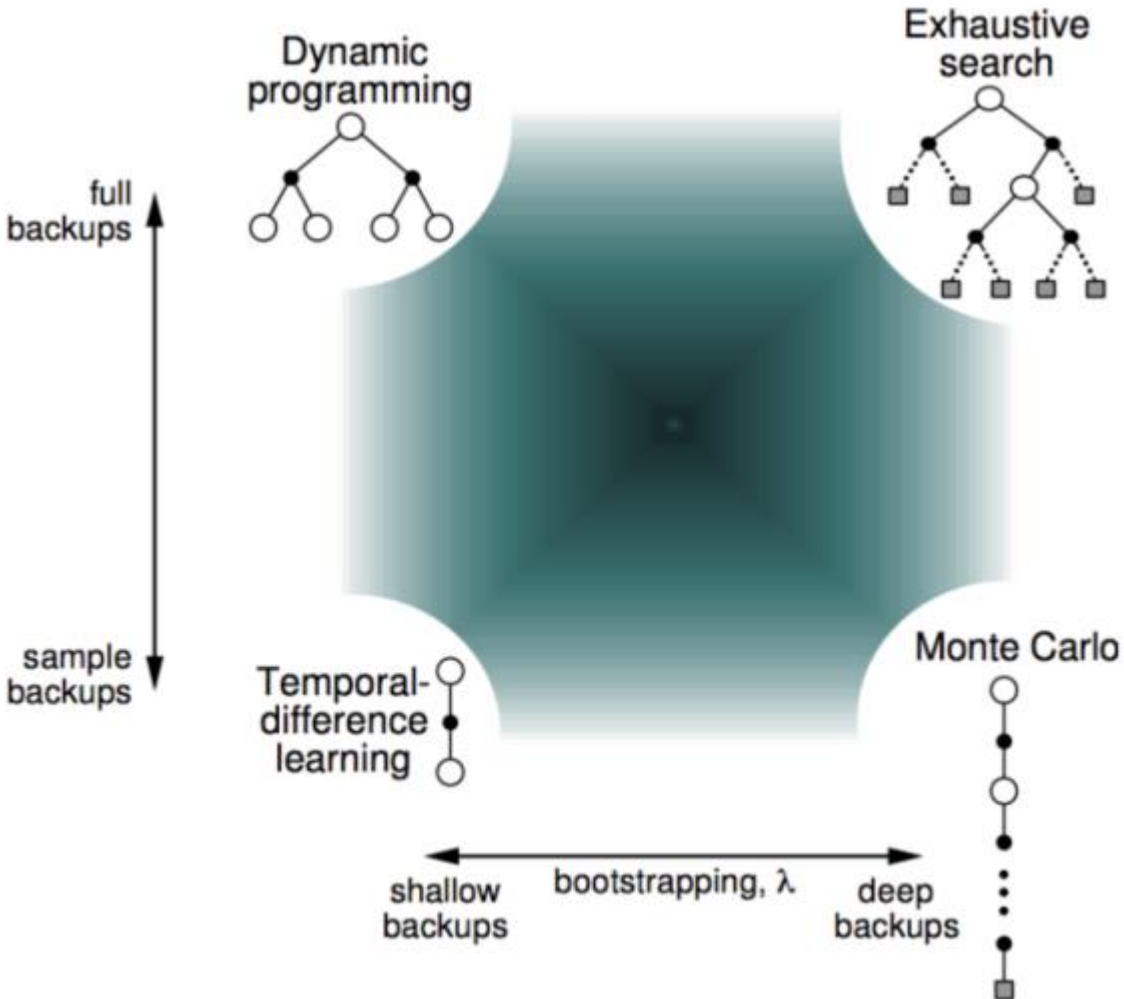
$$V(S_t) \leftarrow \mathbb{E}_{\pi}[R_{t+1} + \gamma V(S_{t+1})]$$



# Bootstrapping and Sampling

- Bootstrapping: update involves an estimate
  - MC does not bootstrap
  - DP bootstraps
  - TD bootstraps
- Sampling: update samples an expectation
  - MC samples
  - DP does not sample
  - TD samples

# Unified View of Reinforcement Learning



# General TD Learning

- Review TD

- Update value  $V(S_t)$  toward estimated return  $R_{t+1} + \gamma V(S_{t+1})$   
$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$
- $R_{t+1} + \gamma V(S_{t+1})$  is called the **TD target**
- For MC learning, the TD target is replaced by  $G_t$   
$$V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$$

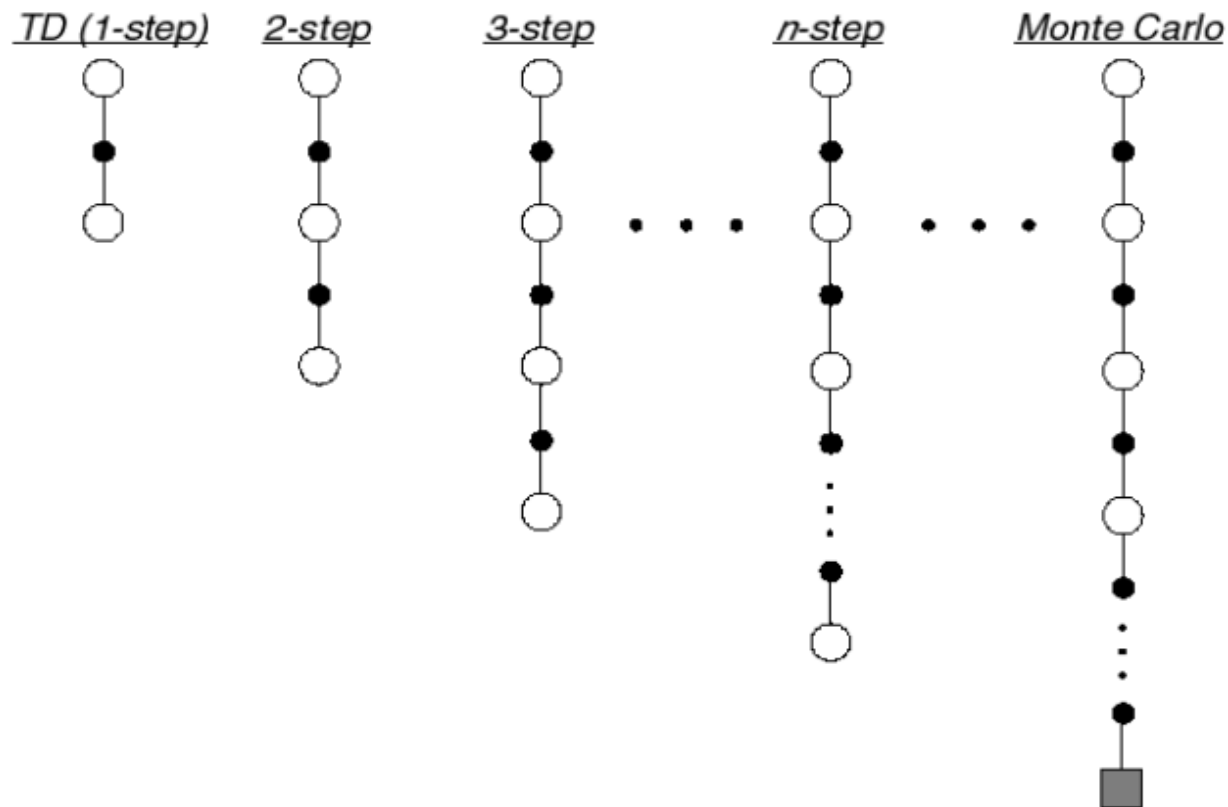
- Question: a more general TD target?
- Investigate TD in a more general manner.
- A typical one: TD( $\lambda$ )





# $n$ -Step Prediction

- Let **TD target** look  $n$  steps into the future



# $n$ -Step Return

- Define the  $n$ -step return

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

- Consider the following  $n$ -step returns for  $n = 1, 2, \infty$

$$n = 1 \quad G_t^{(1)} = R_{t+1} + \gamma V(S_{t+1})$$

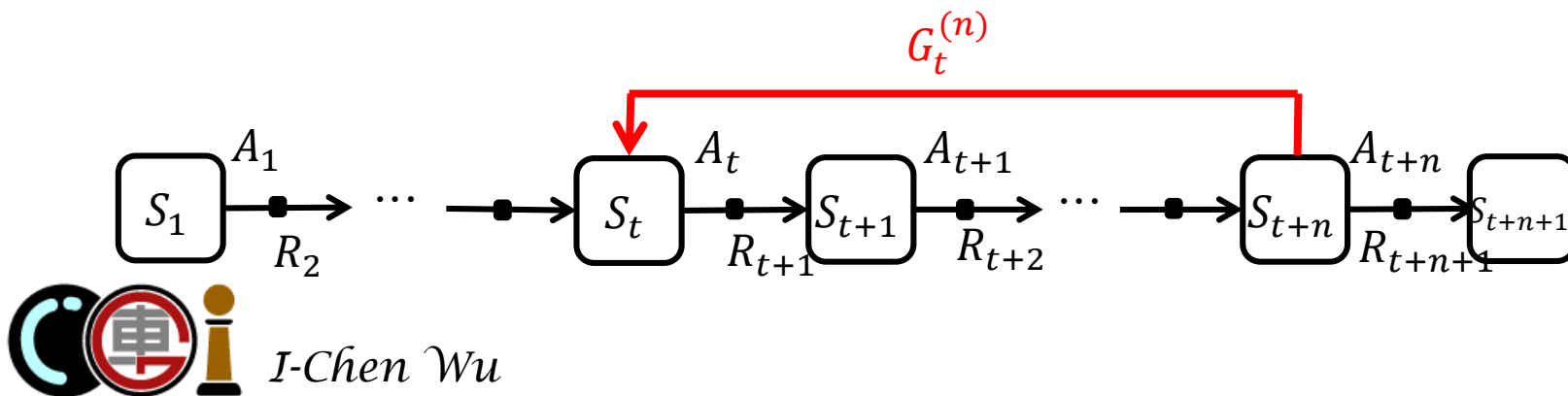
$$n = 2 \quad G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2})$$

$$\vdots$$

$$n = \infty \quad G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T, \text{ if ends at } T.$$

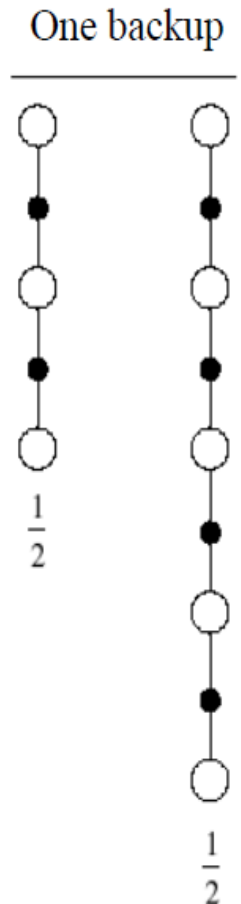
- $n$ -step temporal-difference learning

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t^{(n)} - V(S_t))$$



# Example of Averaging $n$ -Step Returns

- We can average  $n$ -step returns over different  $n$
  - Example:
    - average the 2-step and 4-step returns
- $$\frac{1}{2}G^{(2)} + \frac{1}{2}G^{(4)}$$
- Combines information from two different time-steps
- Next:
    - combine information from all time-steps?



# $\lambda$ -return

- $\lambda$ -return  $G_t^\lambda$ :
  - combines all  $n$ -step returns  $G_t^{(n)}$
- Using weight  $(1 - \lambda) \lambda^{n-1}$

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

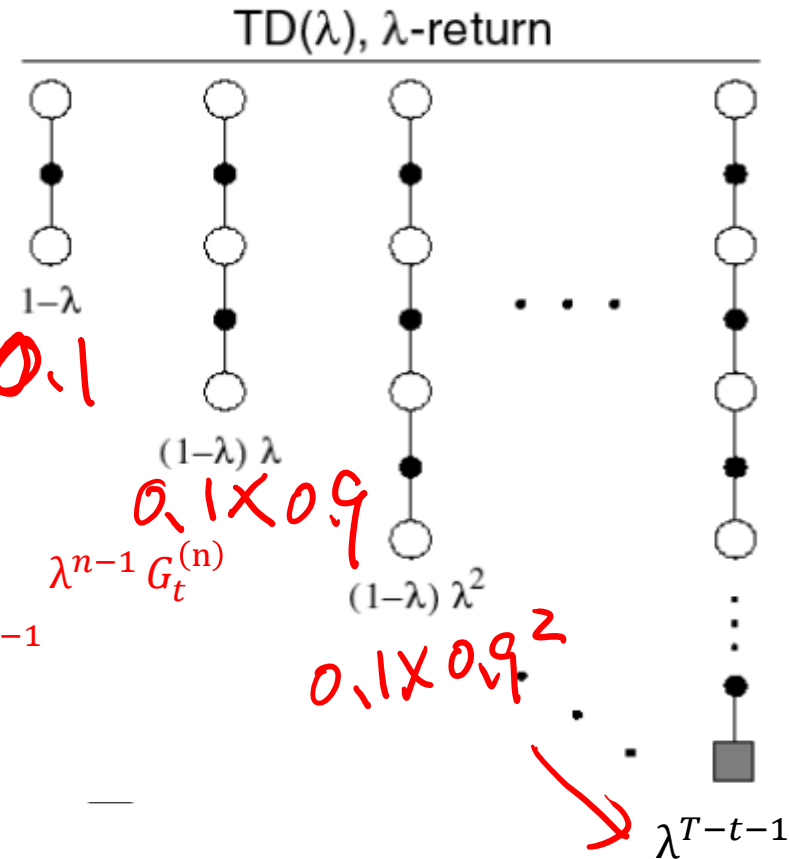
or (in the case of termination)

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_t^{(n)} + (1 - \lambda) \sum_{n=T-t-1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_t^{(n)} + \lambda^{T-t-1} G_t$$

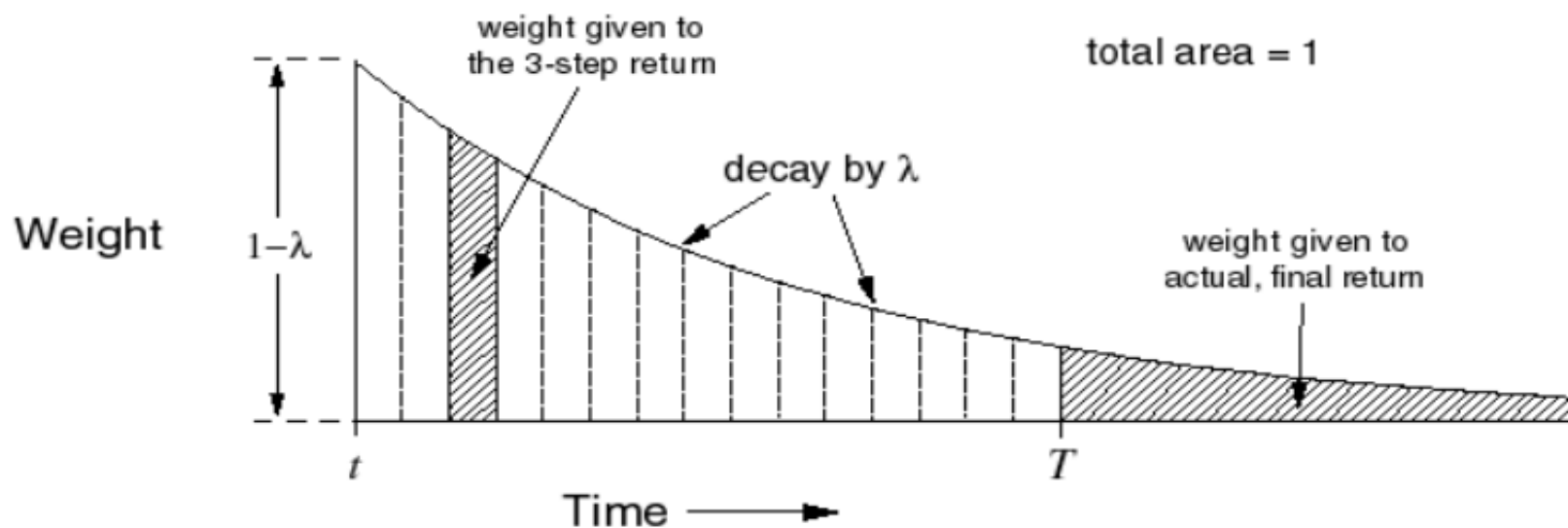
- Forward-view TD( $\lambda$ )

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t^\lambda - V(S_t))$$



# TD( $\lambda$ ) Weighting Function

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$



## TD( $\lambda$ ) and TD(0)

- When  $\lambda = 0$ , only current state is updated

$$V(s) \leftarrow V(s) + \alpha \delta_t$$

- This is exactly equivalent to TD(0) update

## TD( $\lambda$ ) and MC

- When  $\lambda = 0$ , only current state is updated,  $\rightarrow$  TD(0)=TD

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_t^{(n)} + \lambda^{T-t-1} G_t = G_t^{(1)}$$

- This is exactly equivalent to TD target.

- When  $\lambda = 1$ , TD(1) = MC

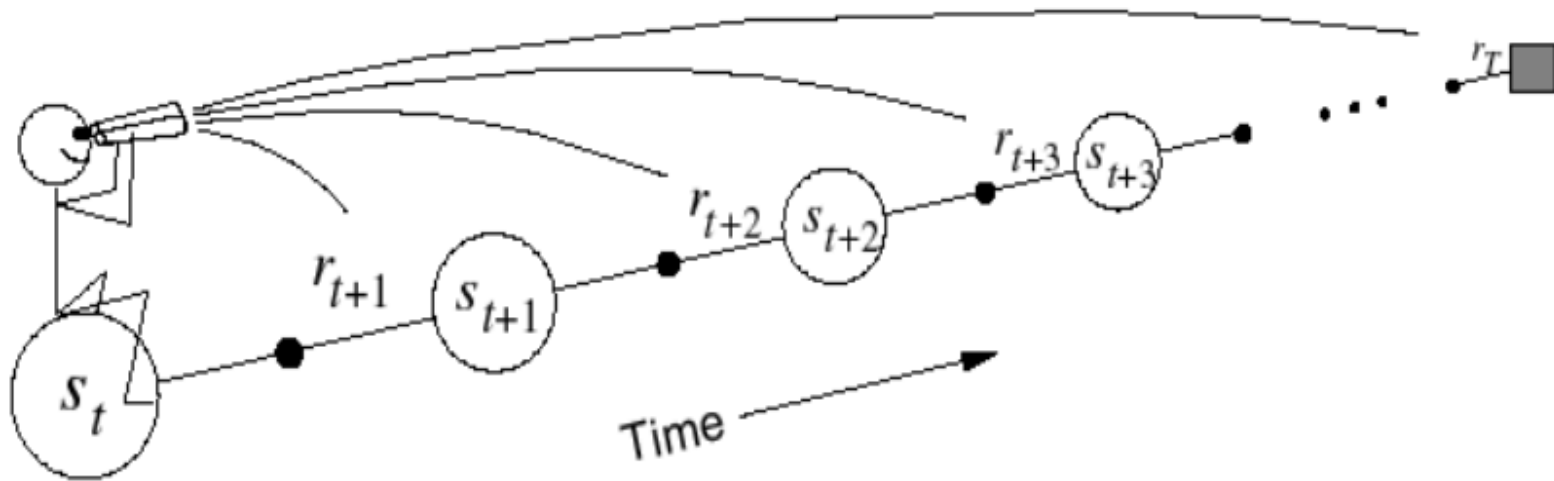
$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_t^{(n)} + \lambda^{T-t-1} G_t = G_t$$

- This is exactly equivalent to MC target.



# Forward-view TD( $\lambda$ )

- Update value function towards the  $\lambda$ -return
- Forward-view looks into the future to compute  $G_t^\lambda$
- Like MC, can only be computed from complete episodes





# Backward View TD( $\lambda$ )

- Forward view provides **theory**
- Backward view provides **mechanism**
  - Update online, every step, from incomplete sequences
- Notes:
  - You may ignore it now.
  - Consider backward (eligible traces) only when you try to implement it.

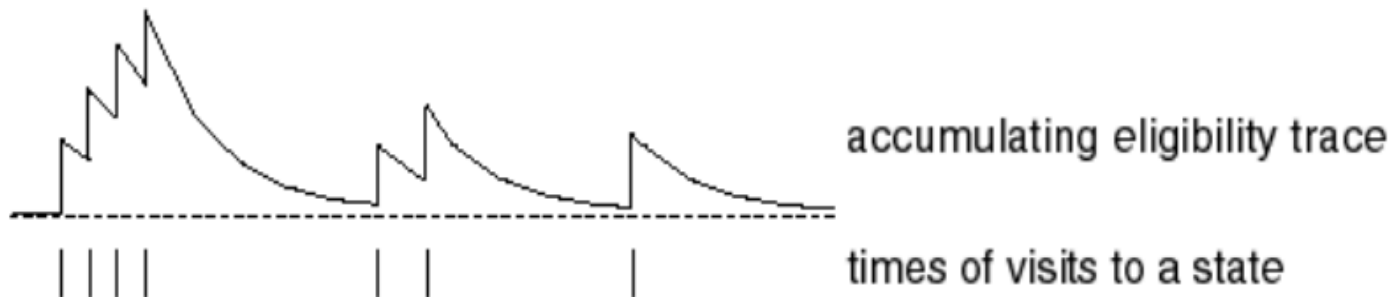
# Eligibility Traces



- Credit assignment problem: did bell or light cause shock?
- **Frequency heuristic**: assign credit to most frequent states
- **Recency heuristic**: assign credit to most recent states
- Eligibility traces combine both heuristics

$$E_0(s) = 0$$

$$E_t(s) = \gamma \lambda E_{t-1}(s) + 1(S_t = s)$$

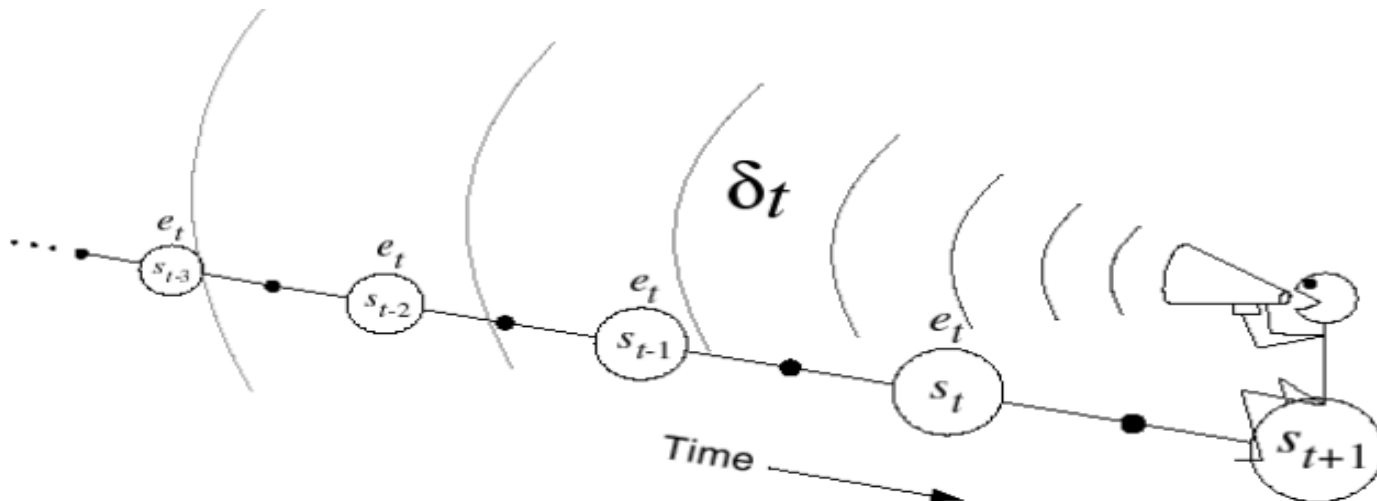


# Backward View TD( $\lambda$ )

- Keep an eligibility trace for every state  $s$
- Update value  $V(s)$  for every state  $s$
- In proportion to TD-error  $\delta_t$  and eligibility trace  $E_t(s)$

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

$$V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$$



# Eligibility Trace

- Explain more in policy-based reinforcement learning for GAE (Generalized Advantage Estimator).
  - (See Page 37, “Advantages and  $TD(\lambda)$ ” in the chapter of policy-based RL.)



---

# $\varepsilon$ -Greedy Exploration



# Example of Greedy Action Selection

- There are two doors in front of you,  
Always apply the greedy action selection:
  - You open the left door and get reward 0  
 $V(left) = 0$
  - You open the right door and get reward +1  
 $V(right) = +1$
  - You open the right door and get reward +3  
 $V(right) = +2$
  - You open the right door and get reward +2  
 $V(right) = +2$
  - $\vdots$
- Are you sure you've chosen the best door?



# $\epsilon$ -Greedy Exploration

- $\epsilon$ -greedy policy:

$$\pi(a|s) = \begin{cases} \epsilon/m + 1 - \epsilon & \text{if } a^* = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q(s, a) \\ \epsilon/m & \text{otherwise} \end{cases}$$

- Exploration

- If you always try the best, you don't explore a real better one.
- With probability  $\epsilon$  choose an action at random
  - ▶ Simplest idea for ensuring continual exploration
- All  $m$  actions are tried with non-zero probability

- Exploitation

- If you always choose at random, you don't exploit the best
- With probability  $1 - \epsilon$  choose the greedy action







# A Lemma in the Previous Proof

(for reference only; can be skipped)

(the sum is a weighted average with nonnegative weights summing to 1, and as such it must be less than or equal to the largest number averaged)

- Lemma: For the previous proof, assume  $\pi(a|s) - \frac{\varepsilon}{m} \geq 0$ .

$$\max_{a \in \mathcal{A}} q_{\pi}(s, a) \geq \sum_{a \in \mathcal{A}} \frac{\pi(a|s) - \frac{\varepsilon}{m}}{1 - \varepsilon} q_{\pi}(s, a)$$

- Proof: Assume all weights  $w_a \geq 0$ , and  $\sum_{a \in \mathcal{A}} w_a = 1$ .

$$\text{Then, } \max_{a \in \mathcal{A}} q_{\pi}(s, a) \geq \sum_{a \in \mathcal{A}} w_a q_{\pi}(s, a)$$

Since weights  $\frac{\pi(a|s) - \frac{\varepsilon}{m}}{1 - \varepsilon} \geq 0$  and their summation = 1,

$$\text{we have } \max_{a \in \mathcal{A}} q_{\pi}(s, a) \geq \sum_{a \in \mathcal{A}} \frac{\pi(a|s) - \frac{\varepsilon}{m}}{1 - \varepsilon} q_{\pi}(s, a)$$



# Key Idea

- Key idea:
  - $\pi(a|s) - \frac{\varepsilon}{m}$  is non-negative, as long as  $\varepsilon$  is monotonically decreasing.
- Example:

Assume  $\varepsilon = 0.4$  and  $m=4$  (4 actions,  $a_1, a_2, a_3, a_4$ ).

  - $\pi(a_1|s) = 0.4$ , and  $q_\pi(a_1|s) = 20$
  - $\pi(a_2|s) = 0.3$ , and  $q_\pi(a_2|s) = 30$  (max in the new policy  $\pi'$ )
  - $\pi(a_3|s) = 0.2$ , and  $q_\pi(a_3|s) = 15$
  - $\pi(a_4|s) = 0.1$ , and  $q_\pi(a_4|s) = 15$
  - Works when  $\varepsilon$  remains the same or drops to a smaller number, say 0.3.



---

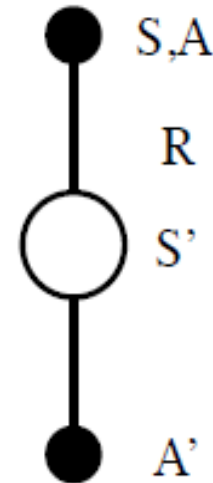
# Q-Learning



# Updating Action-Value Functions with Sarsa

$$Q(S, A) \leftarrow Q(S, A) + \alpha(R + \gamma Q(S', A') - Q(S, A))$$

- Notice: Interesting naming



# Sarsa Algorithm for On-Policy Control

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

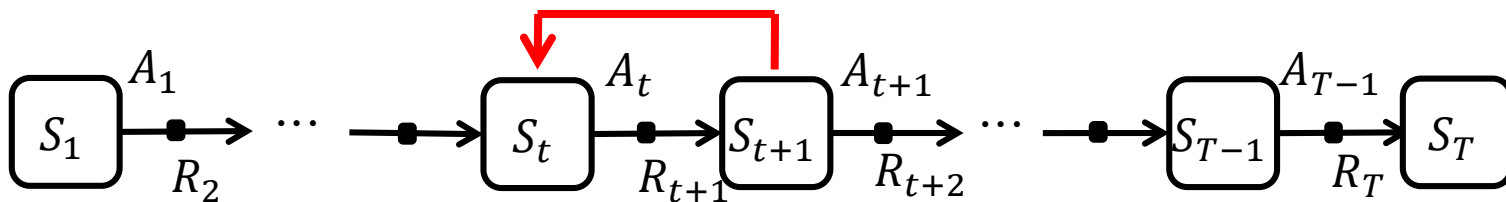
until  $S$  is terminal

- Sarsa converges to the optimal action-value function
- $n$ -step Sarsa – like  $n$ -step return
- Sarsa( $\lambda$ ) – like TD( $\lambda$ )



# Off-Policy Learning

- Evaluate **current policy**  $\pi(a|s)$  to compute  $V_\pi(s)$  or  $q_\pi(s, a)$ , while following **an old policy**  $\mu(a|s)$   
 $\{S_1, A_1, R_2, \dots, S_T\} \sim \mu$
- Why is this important?
  - Learn from observing humans or other agents
  - Re-use experience generated from old policies  $\pi_1, \pi_2, \dots, \pi_{t-1}$
  - Learn about optimal policy while following exploratory policy
  - Learn about multiple policies while following one policy



**Current Policy  $\pi$**



# Importance Sampling

- Estimate the expectation of a different distribution

$$\begin{aligned} & \mathbb{E}_{X \sim P}[f(X)] \\ &= \sum P(X) f(X) \\ &= \sum Q(X) \frac{P(X)}{Q(X)} f(X) \\ &= \mathbb{E}_{X \sim Q} \left[ \frac{P(X)}{Q(X)} f(X) \right] \end{aligned}$$

# Importance Sampling for Off-Policy Monte-Carlo

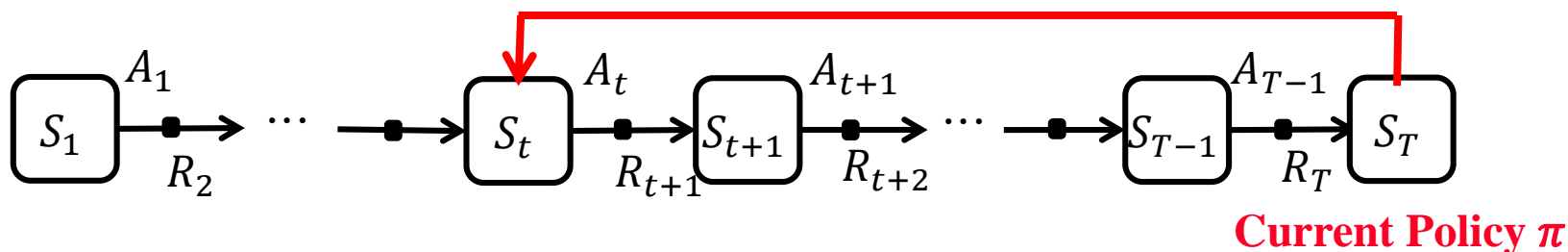
- Use returns generated from  $\mu$  to evaluate  $\pi$
- Weight return  $G_t$  according to similarity between policies
- Multiply importance sampling corrections along whole episode

$$G_t^{\pi/\mu} = \frac{\pi(A_t|S_t)\pi(A_{t+1}|S_{t+1})}{\mu(A_t|S_t)\mu(A_{t+1}|S_{t+1})} \cdots \frac{\pi(A_T|S_T)}{\mu(A_T|S_T)} G_t$$

- Update value towards corrected return

$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t^{\pi/\mu} - V(S_t) \right)$$

- Cannot use if  $\mu$  is zero when  $\pi$  is non-zero
- Importance sampling can dramatically increase variance





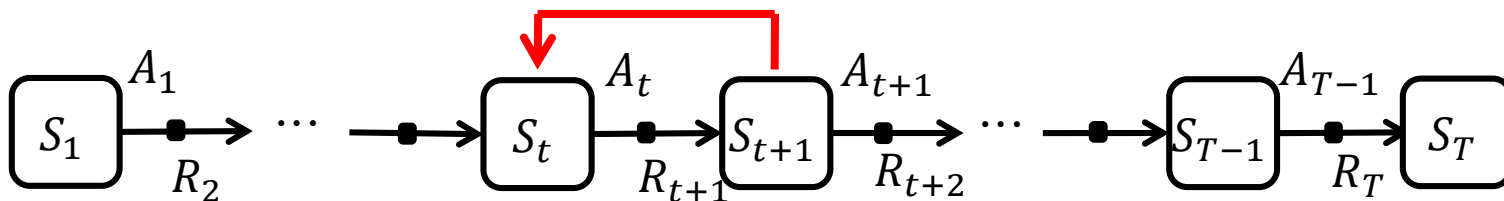
# Importance Sampling for Off-Policy TD

- Use TD targets generated from  $\mu$  to evaluate  $\pi$
- Weight TD target  $R + \gamma V(S')$  by importance sampling
- Only need a single importance sampling correction

$$V(S_t) \leftarrow V(S_t) +$$

$$\alpha \left( \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} (R_{t+1} + \gamma V(S_{t+1})) - V(S_t) \right)$$

- Much lower variance than Monte-Carlo importance sampling (since just one step)
- Policies only need to be similar over a single step



**Current Policy  $\pi$**



# Q-Learning

- We now consider off-policy learning of action-values  $Q(s, a)$
- No importance sampling is required
- Next action is chosen using the old policy
$$A_{t+1} \sim \mu(\cdot | S_{t+1})$$
- But we consider alternative successor action
$$A' \sim \pi(\cdot | S_{t+1})$$
- And update  $Q(S_t, A_t)$  towards value of alternative action
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma Q(S_{t+1}, A') - Q(S_t, A_t))$$

# Off-Policy Control with Q-Learning

- We now allow both old and current policies to improve
- The current policy  $\pi$  is **greedy** w.r.t.  $Q(s, a)$

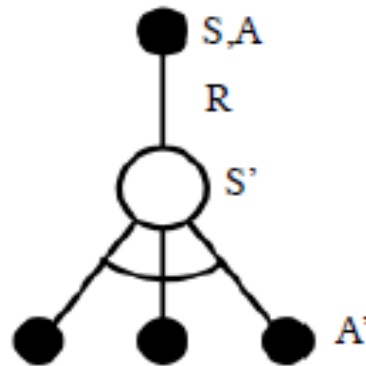
$$\pi(S_{t+1}) = \operatorname{argmax}_{a'} Q(S_{t+1}, a')$$

- The old policy  $\mu$  is e.g.  **$\epsilon$ -greedy** w.r.t.  $Q(s, a)$
- The Q-learning target then simplifies:

$$\begin{aligned} & R_{t+1} + \gamma Q(S_{t+1}, A') \\ &= R_{t+1} + \gamma Q\left(S_{t+1}, \operatorname{argmax}_{a'} Q(S_{t+1}, a')\right) \\ &= R_{t+1} + \max_{a'} \gamma Q(S_{t+1}, a') \end{aligned}$$



# Q-Learning Control Algorithm



- $Q(S, A) \leftarrow Q(S, A) + \alpha \left( R + \gamma \max_{a'} Q(S', a') - Q(S, A) \right)$
- Theorem
  - Q-learning control converges to the optimal action-value function,  $Q(s, a) \rightarrow q_*(s, a)$



# Q-Learning Algorithm for Off-Policy Control

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$   
Repeat (for each episode):  
    Initialize  $S$   
    Repeat (for each step of episode):  
        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)  
        Take action  $A$ , observe  $R, S'$   
         $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$   
         $S \leftarrow S'$   
    until  $S$  is terminal

---

# Function Approximation



# Large-Scale Reinforcement Learning

- Reinforcement learning can be used to solve large problems, e.g.
  - Backgammon:  $10^{20}$  states
  - Computer Go:  $10^{170}$  states
  - Helicopter: continuous state space
- How can we **scale up the model-free methods for prediction and control** from the last two lectures?

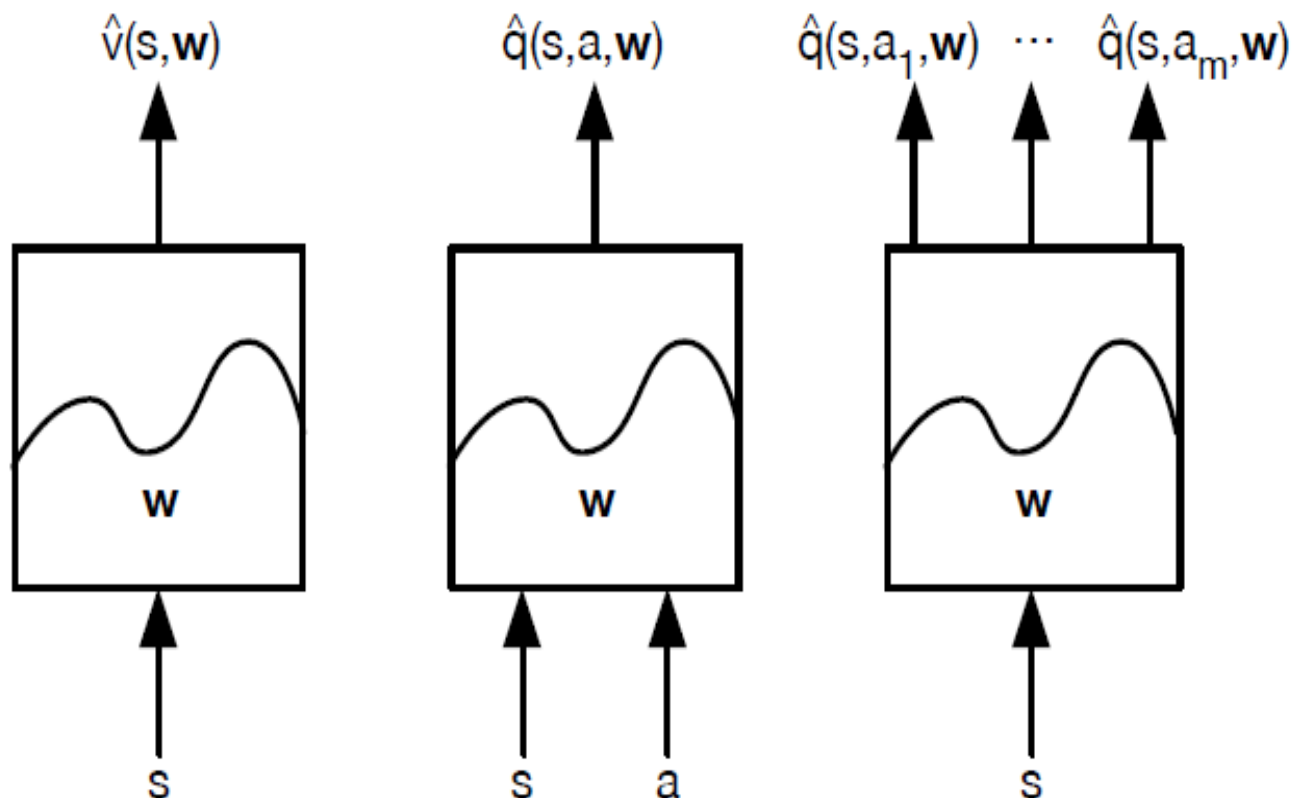
# Value Function Approximation

- So far we have represented value function by a lookup table
  - Every state  $s$  has an entry  $V(s)$
  - Or every state-action pair  $s; a$  has an entry  $Q(s, a)$
- Problem with large MDPs:
  - There are too many states and/or actions to store in memory
  - It is too slow to learn the value of each state individually
- Solution for large MDPs:
  - Estimate value function with function approximation
$$\hat{v}(s, w) \approx v_{\pi}(s)$$
or 
$$\hat{q}(s, a, w) \approx q_{\pi}(s, a)$$
  - Generalize from seen states to unseen states
  - Update parameter  $w$  using MC or TD learning





# Types of Value Function Approximation



# Which Function Approximator?

- There are many function approximators, e.g.
  - Linear combinations of features
  - Neural network
  - Decision tree
  - Nearest neighbour
  - Fourier / wavelet bases
  - ...
- Better to consider **differentiable** function approximators (in red above)
- Furthermore, we require a training method that is suitable for **non-stationary, non-iid data**

# Gradient Descent

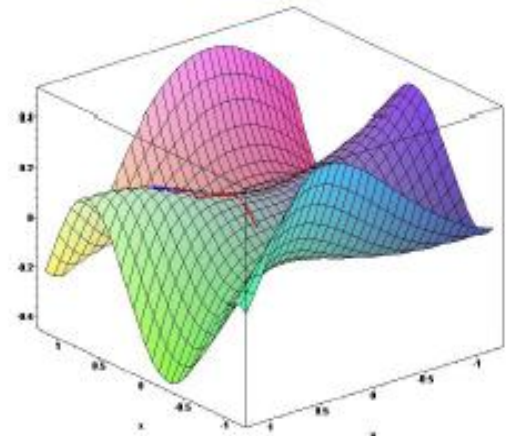
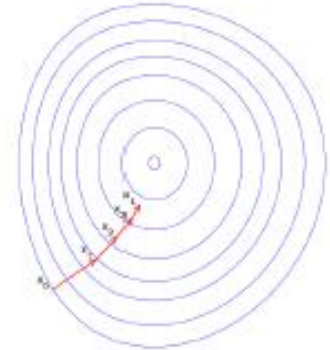
- Let  $J(w)$  be a differentiable function of parameter vector  $w$
- Define the **gradient of  $J(w)$**  to be

$$\nabla_w J(w) = \begin{pmatrix} \frac{\partial J(w)}{\partial w_1} \\ \vdots \\ \frac{\partial J(w)}{\partial w_n} \end{pmatrix}$$

- To find a local minimum of  $J(w)$
- Adjust  $w$  in direction of -ve gradient

$$\Delta w = -\frac{1}{2} \alpha \nabla_w J(w)$$

– where  $\alpha$  is a step-size parameter



# Value Function Approx. By Stochastic Gradient Descent

- Goal: find parameter vector  $w$ 
  - minimizing mean-squared error between approximate value function  $\hat{v}(s, w)$  and true value function  $v_\pi(s)$

$$J(w) = \mathbb{E}_\pi[(v_\pi(S) - \hat{v}(S, w))^2]$$

- Gradient descent **finds a local minimum**

$$\begin{aligned}\Delta w &= -\frac{1}{2} \alpha \nabla_w J(w) \\ &= \alpha \mathbb{E}_\pi[(v_\pi(S) - \hat{v}(S, w)) \nabla_w \hat{v}(S, w)]\end{aligned}$$

- Stochastic gradient descent **samples the gradient**

$$\Delta w = \alpha (v_\pi(S) - \hat{v}(S, w)) \nabla_w \hat{v}(S, w)$$

- Expected update is equal to full gradient update



# Linear Value Function Approximation

- Represent value function by a linear combination of features

$$\hat{v}(S, w) = x(S)^T w = \sum_{j=1}^n x_j(S) w_j$$

- Objective function is quadratic in parameters  $w$

$$J(w) = \mathbb{E}_{\pi}[(v_{\pi}(S) - x(S)^T w)^2]$$

- Stochastic gradient descent converges on global optimum
- Update rule is particularly simple

$$\nabla_w \hat{v}(S, w) = x(S)$$

$$\Delta w = \alpha (v_{\pi}(S) - \hat{v}(S, w)) x(S)$$

- Update = step-size  $\times$  prediction error  $\times$  feature value



# Incremental Prediction Algorithms

- Have assumed true value function  $v_\pi(s)$  given by supervisor
- But in RL there is no supervisor, only rewards
- In practice, we substitute a target for  $v_\pi(s)$

- For MC, the target is the return  $G_t$

$$\Delta w = \alpha (G_t - \hat{v}(S_t, w)) \nabla_w \hat{v}(S_t w)$$

- For TD(0), the target is the TD target  $R_{t+1} + \gamma \hat{v}(S_{t+1}, w)$

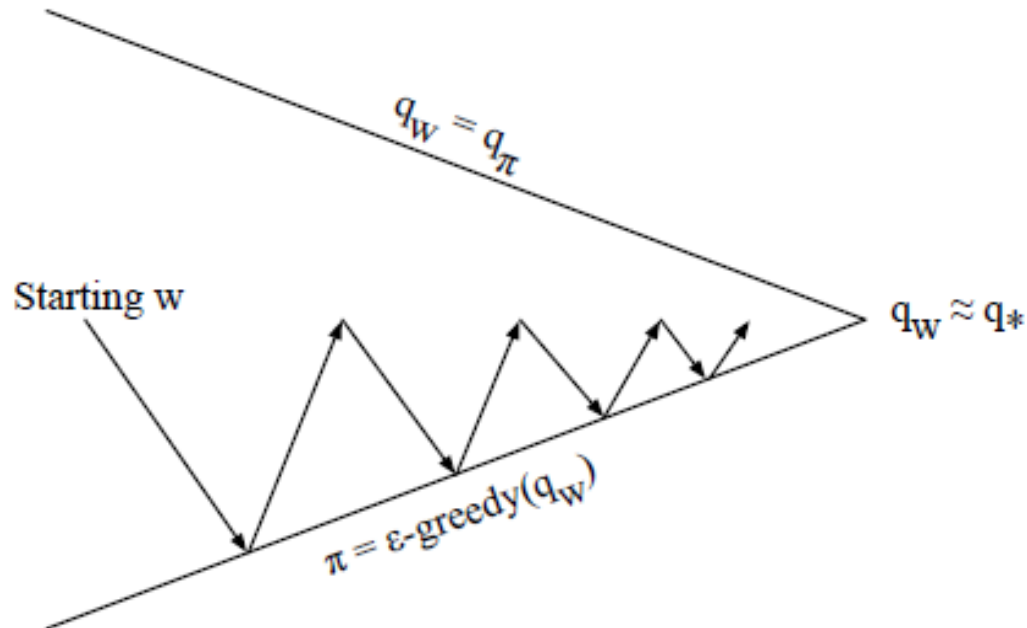
$$\Delta w = \alpha (R_{t+1} + \gamma \hat{v}(S_{t+1}, w) - \hat{v}(S_t, w)) \nabla_w \hat{v}(S_t w)$$

- For TD( $\lambda$ ), the target is the  $\lambda$ -return  $G_t^\lambda$

$$\Delta w = \alpha (G_t^\lambda - \hat{v}(S_t, w)) \nabla_w \hat{v}(S_t w)$$



# Control with Value Function Approximation



- Policy evaluation
  - Approximate policy evaluation,  $\hat{q}(\cdot, \cdot, w) \approx q_\pi$
- Policy improvement
  - $\epsilon$ -greedy policy improvement

# Action-Value Function Approximation

- Approximate the action-value function

$$\hat{q}(S, A, w) \approx q_{\pi}(S, A)$$

- Minimize mean-squared error between approximate action-value function  $\hat{q}(S, A, w)$  and true action-value function  $q_{\pi}(S, A)$

$$J(w) = \mathbb{E}_{\pi}[(q_{\pi}(S, A) - \hat{q}(S, A, w))^2]$$

- Use stochastic gradient descent to find a local minimum

$$-\frac{1}{2} \nabla_w J(w) = (q_{\pi}(S, A) - \hat{q}(S, A, w)) \nabla_w \hat{q}(S, A, w)$$

$$\Delta w = \alpha (q_{\pi}(S, A) - \hat{q}(S, A, w)) \nabla_w \hat{q}(S, A, w)$$





# Incremental Control Algorithms

- Like prediction, we must substitute a target for  $q_\pi(S, A)$

- For MC, the target is the return  $G_t$

$$\Delta w = \alpha (G_t + \hat{q}(S_t, A_t, w)) \nabla_w \hat{q}(S_t, A_t, w)$$

- For TD(0), the target is the TD target  $R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$

$$\Delta w = \alpha (R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, w) - \hat{q}(S_t, A_t, w)) \nabla_w \hat{q}(S_t, A_t, w)$$

- For forward-view TD( $\lambda$ ), target is the action-value  $\lambda$ -return

$$\Delta w = \alpha (q_t^\lambda - \hat{q}(S_t, A_t, w)) \nabla_w \hat{q}(S_t, A_t, w)$$

- For backward-view TD( $\lambda$ ), equivalent update is

$$\delta_t = R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, w) - \hat{q}(S_t, A_t, w)$$

$$E_t = \gamma \lambda E_{t-1} + \nabla_w \hat{q}(S_t, A_t, w)$$

$$\Delta w = \alpha \delta_t E_t$$



# Batch Reinforcement Learning

- Gradient descent is simple and appealing
- But it is **not sample efficient**
- **Batch methods seek to find the best fitting value function**
- Given the agent's experience (“training data”)

# Least Squares Prediction

- Given value function approximation  $\hat{v}(s, w) \approx v_\pi(s)$
- And experience  $D$  consisting of  $\langle \text{state}, \text{value} \rangle$  pairs
$$D = \{\langle s_1, v_1^\pi \rangle, \langle s_2, v_2^\pi \rangle, \dots, \langle s_T, v_T^\pi \rangle\}$$
- Which parameters  $w$  give the best fitting value fn  $\hat{v}(s, w)$ ?
- **Least squares algorithms** find parameter vector  $w$  minimizing sum-squared error between  $\hat{v}(s_t, w)$  and target values  $v_t^\pi$ ,

$$\begin{aligned} LS(w) &= \sum_{t=1}^T (v_t^\pi - \hat{v}(s_t, w))^2 \\ &= \mathbb{E}_D [(v^\pi - \hat{v}(s, w))^2] \end{aligned}$$

# Stochastic Gradient Descent with Experience Replay

- Given experience consisting of  $\langle \text{state}, \text{value} \rangle$  pairs

$$D = \{ \langle s_1, v_1^\pi \rangle, \langle s_2, v_2^\pi \rangle, \dots \langle s_T, v_T^\pi \rangle \}$$

- Repeat:

- Sample state, value from experience

$$\langle s, v^\pi \rangle \sim D$$

- Apply stochastic gradient descent update

$$\Delta w = \alpha (v^\pi - \hat{v}(s, w)) \nabla_w \hat{v}(s, w)$$

- Converges to least squares solution

$$w^\pi = \underset{w}{\operatorname{argmin}} LS(w)$$

- Similar for action value function  $q^\pi$

