CS395T: Foundations of Machine Learning for Systems Researchers

Fall 2025

## **Lecture 7:**

Model-free (Sampling) Methods for MDPs

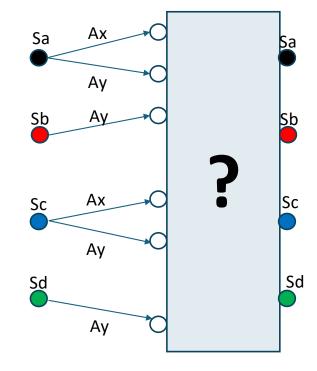


## Problem

In practice, we usually have incomplete knowledge of environment

- Rewards?
- Probabilities?

How do we do policy optimization when we have incomplete knowledge of environment?



## Model-free methods

Based on sampling: determine optimal policy from estimates of rewards and probabilities

## Issues

- What are samples?
- How do we learn from a sample?
- When do we learn from samples?
- How to update policy when you get new samples?

Goal is to determine optimal policy and not to learn full MDP, which is usually impossible



## Organization

## Background:

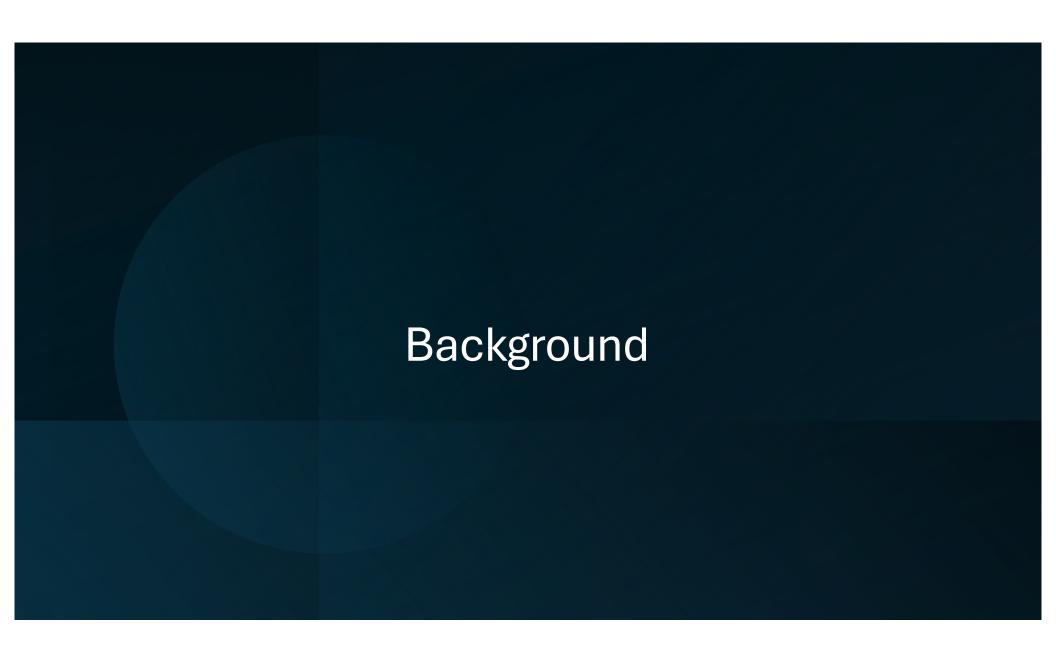
 Incremental & exponential recency-weighted averaging of sequences

## Offline method

- Collect samples and post-process them to find V/Q estimates
- Like batching in NN training

## Online methods

- Boot-strapping methods: start with initial approximations for V/Q values and update when you get sample
  - Temporal-difference (TD) methods: TD(n),  $TD(\lambda)$ , SARSA, Q-learning
- Non-bootstrapping method: no initial approximation to V/Q values needed
  - Monte Carlo (RL terminology)



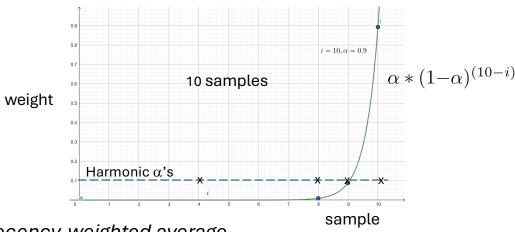
## Incremental averaging of sequence of values

- Problem: given
  - Sequence of (possibly noisy) measurements of some quantity:  $V_1, V_2, V_3, ...$
  - Maintain a running average of measurements:  $\hat{V}_1, \hat{V}_2, \hat{V}_3, ...$

$$\hat{V}_{1} = V_{1} 
\hat{V}_{i} = \frac{V_{1} + \dots + V_{i}}{i} \quad (i \ge 1) 
= \frac{V_{1} + \dots + V_{i-1}}{i-1} * \frac{i-1}{i} + \frac{V_{i}}{i} 
= (1 - \frac{1}{i})\hat{V}_{i-1} + \frac{1}{i} * V_{i}$$

Target 
$$\hat{V}_i = \hat{V}_{i-1} + \alpha_i (\hat{V}_i - \hat{V}_{i-1}) \quad \text{where } \alpha_i = \frac{1}{i}$$
 Harmonic  $\alpha$ 's

# Non-stationary problems



- Exponential recency-weighted average
  - For non-stationary problems, better to give more weight to more recent measurements
  - · Maintain a weighted average of measurements

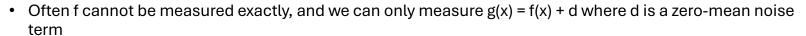
$$\hat{V}_{1} = V_{1} 
\hat{V}_{i} = \hat{V}_{i-1} + \alpha(V_{i} - \hat{V}_{i-1}) \text{ where } 0 \le \alpha \le 1 \text{ is constant} 
\hat{V}_{i} = \alpha V_{i} + \alpha(1 - \alpha)V_{i-1} + \alpha(1 - \alpha)^{2}V_{i-2} + \dots + \alpha(1 - \alpha)^{i-2}V_{2} + (1 - \alpha)^{i-1}V_{1}$$

 $\alpha \to 1$  gives more weight to recent measurements

 $\alpha \to 0$  gives more weight to older measurements

# Stochastic Approximation Theory

- Given: continuous function f(x), non-decreasing in interval (l,h), and f(l) < 0 and f(h) > 0
- Intermediate value theorem:  $\exists x^*: l < x^* < h$  such that  $f(x^*) = 0$
- Iterative method for finding x\*
  - $x_1$  = any point in interval (l,h)
  - if  $f(x_1) < 0$  try point to right of  $x_1$  otherwise try point to left of  $x_1$ 
    - Iterative scheme:  $x_i = x_{i-1} \alpha^* f(x_{i-1})$  where  $\alpha$  is some constant > 0

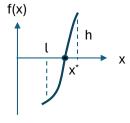


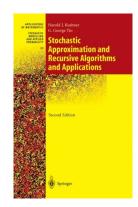
- One solution: measure  $g(x_i)$  many times and take average
  - Problem: you may spend a lot of effort in estimating  $f(x_i)$  even if  $x_i$  is far from  $x^*$
- Robbins-Monro[1951]: use iterative scheme with adaptive  $\alpha_i > 0$

$$x_i = x_{i-1} - \alpha_i * g(x_{i-1})$$
 where  $\sum_{i=1}^{\infty} \alpha_i = \infty$  and  $\sum_{i=1}^{\infty} \alpha_i^2 < \infty$   $x_i$  converges to  $x^*$ 

• Stochastic gradient-descent: f is the derivative of loss function



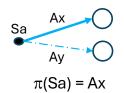




# Sampling Methods



## High-level idea



Focus on *policy evaluation* for policy  $\pi$ ,  $H = \infty$ : need to solve linear system

$$V^{\pi}(s) = \sum_{s'} P(s, \pi(s), s') * (R(s, \pi(s), s') + \gamma * V^{\pi}(s'))$$

Since  $\pi$  is fixed, simplify notation by dropping  $\pi$  from equations

$$V(s) = \sum_{s'} P(s, s') * (R(s, s') + \gamma * V(s'))$$

We do not know P and R, but assume

 $\hat{P}$ : estimate for P and row-stochastic matrix

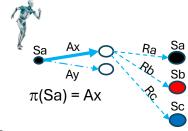
 $\hat{R}$ : estimate for R

Compute Û, an estimate for state valuations, by solving linear system

$$\hat{V}(s) = \sum_{s'} \hat{P}(s, s') * (\hat{R}(s, s') + \gamma * \hat{V}(s'))$$

If  $\hat{P} \approx P$  and  $\hat{R} \approx R$ , then  $\hat{V} \approx V$ 

## First attempt: offline method



### Collect a multiset of samples

- At state Sa, agent takes action  $\pi(Sa)$  and sees what happens
  - Observation SARS = <StartState, Action, Observed Reward, Observed State>
- Repeat to obtain multiset of observations starting at various states
  - $O = [SARS_1, SARS_2, SARS_3, \dots, SARS_n]$

## Estimate V by offline processing of observations

- n(Sx) = number of times agent started in state Sx in O (assume > 0)
- n(Sx,Sy) = number of times agent ended up in state Sy when it started in state Sx
- $\hat{P}(Sx,Sy) = \frac{n(Sx,Sy)}{n(Sx)}$  ( $\hat{P}$  is a row-stochastic matrix)
- $\hat{R}(Sx,Sy) = observed reward in any Sx \rightarrow Sy sample$
- Solve linear system to find  $\hat{V}$

$$\hat{V}(s) = \frac{1}{n(s)} \sum_{s'} n(s, s') * (\hat{R}(s, s') + \gamma * \hat{V}(s'))$$

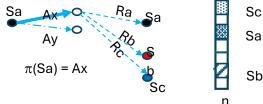
## Drawbacks of offline method



- 1. We do not learn valuations until all samples have been collected.
  - Online methods: pipeline sample collection and valuation estimation
  - Exploitation: use valuations to guide where to sample
- 2. Inefficient to jump around between starting states while collecting samples
  - Solution: agent should follow paths in the MDP graph while collecting samples rather than jumping around
  - Episodes

For now, collect samples, store on disk and process them in streaming fashion iteratively until convergence.
Like "batching" in NN training

## TD(0): simplest online method



Recall: offline processing

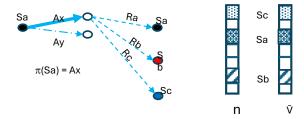
$$\hat{V}(s) = \frac{1}{n(s)} \sum_{s'} n(s, s') * (\hat{R}(s, s') + \gamma * \hat{V}(s'))$$

- Online updates: maintain two arrays while processing samples
  - n(s) = number of times state s has been visited (initialized to 0)
  - $\tilde{V}(s)$ = current estimate of V (initialized arbitrarily)|
- New sample <Sa, Ax, R(Sa,S'), S'> comes in

$$\tilde{V}(Sa) \leftarrow \frac{n(Sa)*\tilde{V}(Sa) + (R(Sa,S') + \gamma*\tilde{V}(S'))}{n(Sa) + 1} \text{ which is equivalent to } \\ \tilde{V}(Sa) \leftarrow \tilde{V}(Sa) + \frac{1}{\alpha}(\underline{R(Sa,S') + \gamma*V(S') - V(Sa)}) \text{ where } \alpha = n(Sa) + 1 \\ n(Sa) \leftarrow n(Sa) + 1$$

Intuition: average over multiset of samples = weighted average over sample set, weighted by frequency

# TD(0) program



```
Array n:1..|S| = 0; //number of visits to states Array \tilde{V}:1..|S| = \text{arbitrary}; //V estimates while (not converged) do for each sample \langle \text{Sa}, \pi(\text{Sa}), \text{R}(\text{Sa}, \text{S}'), \text{S}' \rangle do { increment n(\text{Sa}); \tilde{V}(\text{Sa}) \leftarrow \tilde{V}(\text{Sa}) + \frac{1}{n(Sa)} (\text{R}(\text{Sa}, \text{S}') + \gamma * \tilde{V}(\text{S}') - \tilde{V}(\text{Sa})); }
```

Tabular method: maintain table to map states to values

## Comments

Convergence (Balachander, Chatterjee et al.): Given a sequence of one-hop samples O, the  $\tilde{V}$  values computed by the TD(0) program will converge to the  $\hat{V}$  values in the fixpoint equation

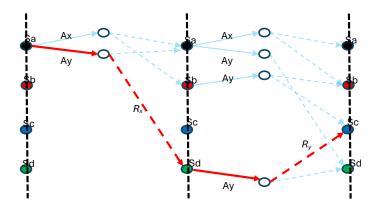
Caveat: need to iterate over sample sequence O many times like in batching

Convergence:  $\tilde{V}$  values computed by TD(0) program will converge to  $V^{\pi}$  if each state is visited unboundedly often in O (called *exploring starts*).

In practice: throw away sample after it is processed and regenerate later

Boot-strapping method:  $\tilde{V}$  is initialized arbitrarily Biased estimator until convergence

# TD(1)

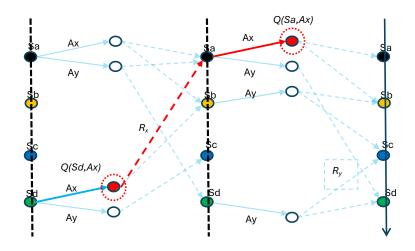


- TD(1) two-hop sample: <S<sub>0</sub>,A<sub>0</sub>,R<sub>1</sub>,S<sub>1</sub>,A<sub>1</sub>,R<sub>2</sub>,S<sub>2</sub>> using policy  $\pi$
- · Update valuations by computing total discounted reward

$$V(Sa) \leftarrow V(Sa) + \alpha[Rx + \gamma * Ry + \gamma^2 V(Sc) - V(Sa)] \text{ (where } \alpha = n(Sa))$$

- Generalize to TD(n): paths with (n-1) hops
- Convergence arguments are similar to TD(0) case

## SARSA



- Maintain table of Q-values rather than V values, using policy  $\boldsymbol{\pi}$ 

$$Q(Sd, Ax) \leftarrow Q(Sd, Ax) + \alpha * [Rx + \gamma * Q(Sa, \pi(Sa)) - Q(Sd, Ax)]$$

• Generalize to multiple hops

# On-policy and off-policy methods

## **Terminology**

- · Target policy: policy you are optimizing
- · Behavior policy: policy that generates actions

## On-policy methods

- Target policy = Behavior policy
- "Learn about the policy you are following" (alternate policy evaluation and improvement)
- Can be less efficient than off-policy methods
- Example: TD(n), SARSA

## Off-policy methods

- Behavior policy ≠ Target policy
- "Learn about one policy while behaving according to different policy" (fuse policy evaluation & improvement)
- · May find optimal policy faster than on-policy methods
- · Examples: Q-learning

Good informal explanation

## **Exploratory policies**

**Pure exploration**: pick an action at random. Does not exploit Q-value estimates.

 $\epsilon$ -greedy: use exploration  $\epsilon$  (~5-10%) fraction of time, and exploitation rest of time. For exploration, choose action at random. Most popular method.

**Boltzmann exploration**: like  $\epsilon$ -greedy but uses softmax over current Q-value estimates to select action for exploration

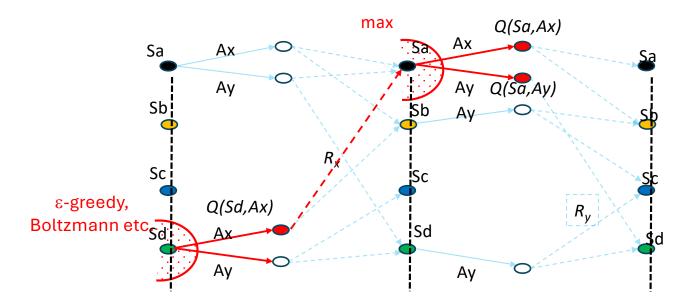
**Pure exploitation**: pick action with highest Q-value. May get stuck in local minimum.

Policy network: DNN maps state to distribution over actions, and sampler picks action

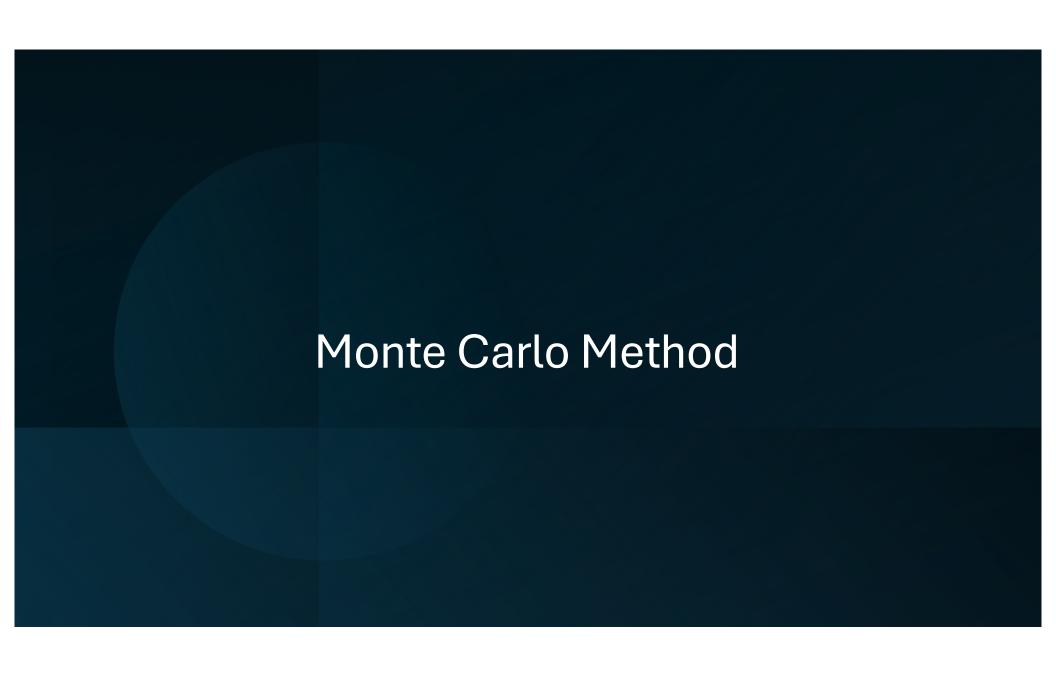
**Exploration-exploitation tradeoff** 

Good explanation of off-policy methods

# Q-learning



- Off-policy method: needs exploratory policy
- Compute Q-values like SARSA but optimize over actions like value iteration  $Q(Sd,Ax) \ \leftarrow \ Q(Sd,Ax) + \alpha * [Rx + \gamma * max_{As}Q(Sa,As) Q(Sd,Ax)]$
- Generalize to multiple hops
- Convergence of O-learning



# **Expectations over Trajectories**

Unrolled MDP DAG for fixed  $\pi$ , rooted at start state S0, H = T

Each edge is labeled with probability and reward Sum of probabilities on outgoing edges of node = 1

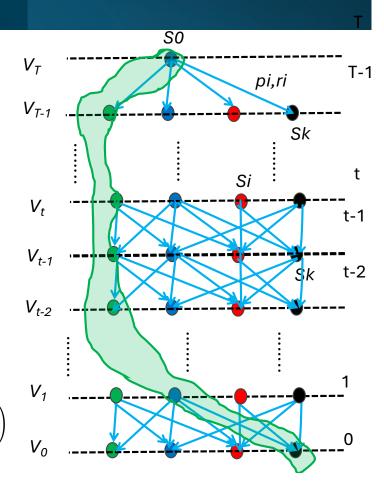
*Trajectory*:  $\tau \sim \pi$  path from root to a leaf

Probability of trajectory  $\tau$ :  $\mathbb{P}(\tau)$  = product of edge probabilities

Reward of trajectory  $\tau$ : R( $\tau$ ) = sum of rewards on trajectory

Sum of trajectory probabilities = 1

Meaningful to talk about  $\mathbb{E}_{\tau \sim \pi} \Big[ \text{property of } \tau \Big] \bigg( = \sum_{\tau \sim \pi} \mathbb{P}(\tau) (\text{property of } \tau) \bigg)$ 

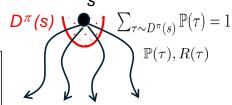


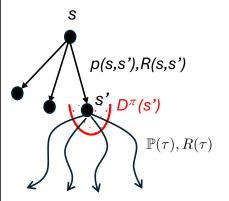
# V/Q Valuations as Expectations over Trajectories: Continuous tasks

$$V^{\pi}(s) = \mathbb{E}_{\tau \sim D^{\pi}(s)} \Big[ R(\tau) \Big]$$
 (where  $D^{\pi}(s)$  is set of trajectories starting at s)

Proof for continuous tasks (episodic task proof is similar)

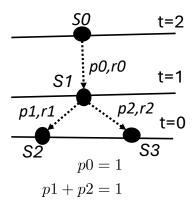
$$\begin{split} V^{\pi}(s) &= \sum_{s'} p(s,s') [ \ R(s,s') + \gamma \ V^{\pi}(s') \ ] \quad \text{(fixpoint equation)} \\ &= \sum_{s'} p(s,s') [ \ R(s,s') + \gamma \ \sum_{\tau \sim D^{\pi}(s')} \mathbb{P}(\tau) * R(\tau) \ ] \quad \text{(assumption about } V^{\pi}(s') \text{)} \\ &= \sum_{s'} p(s,s') [ \ \sum_{\tau \sim D^{\pi}(s')} \mathbb{P}(\tau) * R(s,s') + \gamma \ \sum_{\tau \sim D^{\pi}(s')} \mathbb{P}(\tau) * R(\tau) \ ] \quad \text{(because } \left[ \sum_{\tau \sim D^{\pi}(s')} \mathbb{P}(\tau) = 1 \right) \\ &= \sum_{s'} p(s,s') \sum_{\tau \sim D^{\pi}(s')} \mathbb{P}(\tau) * \left[ (R(s,s') + \gamma * R(\tau)) \right] \\ &= \sum_{s'} \sum_{\tau \sim D^{\pi}(s)} p(s,s') * \mathbb{P}(\tau) * \left[ (R(s,s') + \gamma * R(\tau)) \right] \\ &= \sum_{\tau \sim D^{\pi}(s)} \mathbb{P}(\tau) * (R(\tau)) \end{split}$$





## Example

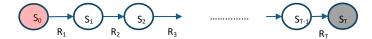
Bellman: 
$$V(S0) = \boxed{p0*(r0 + (p1*r1 + p2*r2))}$$
 
$$= p0*r0 + p0*p1*r1 + p0*p2*r2$$
 Trajectories: 
$$V(S0) = \boxed{p0*p1*(r0+r1) + p0*p2*(r0+r2)}$$
 
$$= p0*p1*r0 + p0*p1*r1 + p0*p2*r0 + p0*p2*r2$$
 
$$= p0*r0*\underbrace{(p1+p2)}_{=1} + p0*p1*r1 + p0*p2*r2$$
 
$$= p0*r0 + p0*p1*r1 + p0*p2*r2$$



Unrolled MDP for H = 2, policy  $\pi$ 

Intuition: rooted DAG  $\rightarrow$  set of paths starting at root Contribution of each edge to expected reward at S0 is distributed over trajectories that contain it

## Episode



From a mathematical perspective, multiset of samples can be obtained by jumping randomly between states

Practical view: more economical to follow **paths in the state space**, recording states and rewards, optionally updating V/Q values on the fly

### Terminal/absorbing states

- States with no outgoing edges in MDP graph
- (E.g.) win/lose states in two-person games

**Episode:** <S<sub>0</sub>,A<sub>0</sub>,R<sub>1</sub>,S<sub>1</sub>,A<sub>1</sub>,R<sub>2</sub>,S<sub>2</sub>,....,S<sub>T-1</sub>,A<sub>T-1</sub>,R<sub>T</sub>,S<sub>T</sub>>

- Sequence of state transitions that form a path in the MDP graph
- S<sub>T</sub> is terminal state

## Monte Carlo Method



Problem Terminal state

- MDP with absorbing/terminal states w/fixed valuations
- Policy is fixed:  $\pi$

#### Monte Carlo method

- Sample multiset of episodes
  - At each state, sum discounted rewards from all samples starting at that state
  - · At the end, divide sum by number of episodes starting at that state
- · Intuition: expected number of times path is sampled is proportional to its probability
- Non-boot-strapping method: we are propagating updates to state valuations back from terminal states with fixed valuations

## **Details**

- For given episode, update valuations of all intermediate nodes on path
  - Implementation: propagate (discounted) rewards backwards along path
- Repeated occurrences of node on path
  - First-visit MC vs. every-visit MC

# Monte Carlo (no discount): Barto & Sutton

```
Algorithm 1: First-Visit MC Prediction
 Input: policy \pi, positive integer num\_episodes
 Output: value function V \approx v_{\pi}, if num\_episodes is large enough)
 Initialize N(s) = 0 for all s \in \mathcal{S}
 Initialize Returns(s) = 0 for all s \in \mathcal{S}
 for episode\ e \leftarrow 1 to e \leftarrow num\_episode\ do
      Generate, using \pi, an episode S_0, A_0, R_1, S_1, A_1, R_2, \ldots, S_{T-1}, A_{T-1}, R_T
      G \leftarrow 0
      for time step t = T - 1 to t = 0 (of the episode e) do
          G \leftarrow G + R_{t+1}
          if state S_t is not in the sequence S_0, S_1, \ldots, S_{t-1} then
                Returns(S_t) \leftarrow Returns(S_t) + G_t
                N(S_t) \leftarrow N(S_t) + 1
      end
  \quad \mathbf{end} \quad
 V(s) \leftarrow \frac{\text{Returns}(s)}{N(s)} \text{ for all } s \in \mathcal{S}
\underline{\mathbf{retur}}n V
```

# $\begin{aligned} & \textbf{Algorithm 2: Every-Visit MC Prediction} \\ & \textbf{Input: policy $\pi$, positive integer $num\_episodes$} \\ & \textbf{Output: value function $V$ $(\approx v_{\pi}$, if $num\_episodes$ is large enough)} \\ & \textbf{Initialize $N(s) = 0$ for all $s \in \mathcal{S}$} \\ & \textbf{Initialize Returns}(s) = 0 \text{ for all } s \in \mathcal{S} \\ & \textbf{for $episode $e \leftarrow 1$ $to $e \leftarrow num\_episodes$ $\textbf{do}$} \\ & \textbf{Generate, using $\pi$, an episode $S_0, A_0, R_1, S_1, A_1, R_2 \dots, S_{T-1}, A_{T-1}, R_T$} \\ & \textbf{G} \leftarrow 0 \\ & \textbf{for $time $step $t = T - 1$ $to $t = 0$ (of the $episode $e$) $\textbf{do}$} \\ & \textbf{Returns}(S_t) \leftarrow \text{Returns}(S_t) + G_t \\ & \textbf{N}(S_t) \leftarrow \textbf{N}(S_t) + 1 \\ & \textbf{end} \end{aligned}$

## Pros and cons of MC vs. TD

## TD can learn before knowing terminal state in episode

- TD can learn after every step of episode
- MC must wait till the terminal state is known

## TD can learn even if there is no terminal state in episode

- TD can learn even in continuing environments where there are no terminal states
- MC only works for episodic (terminating) environments

#### Bias/variance tradeoff

- Bias
  - MC+ : Returns from terminating paths: unbiased estimator of  $V^{\pi}$
  - TD-: TD targets  $R_{t+1} + \gamma^* V_i^{\pi}(S_t)$ : biased estimator (until convergence)
- Variance
  - MC-: episodes may have long and variable length paths, so high variance
  - TD+: short paths, so lower variance

# Lecture 4: Model-Free Prediction Lecture 4: Model-Free Prediction

# Estimated value 0.4 - 0.2 - 0.2 - 0.2 - 0.8 - 0.6 - 0.4 - 0.2 - 0.

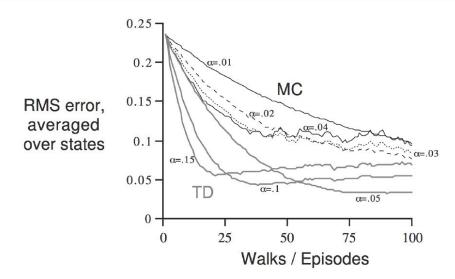
В

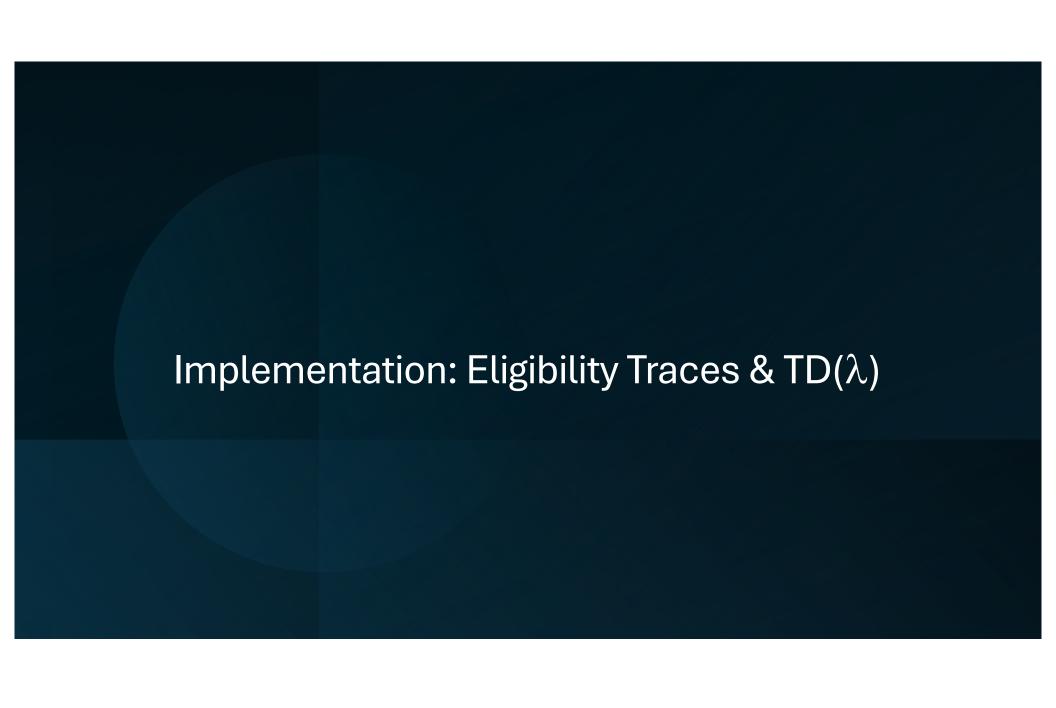
Ċ State E

D

# Convergence







## TD(0) Pseudocode with Episodes (Barto & Sutton)

```
Algorithm 1 Tabular TD(0) for estimating v_{\pi}

Input: Policy \pi to be evaluated Parameters: Learning rate \alpha \in (0, 1]

1: for each episode: do

2: Initialize S

3: while S is not terminal: do

4: Take action A given by \pi(a|S)

5: Observe R, S'

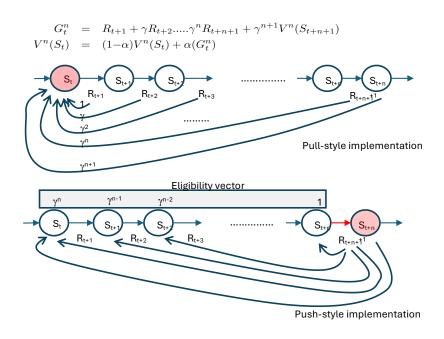
6: Update V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]

7: S \leftarrow S'

8: end while

9: end for
```

# Implementing TD(n): push vs. pull



## Detail: updates actually performed using TD(0) Errors

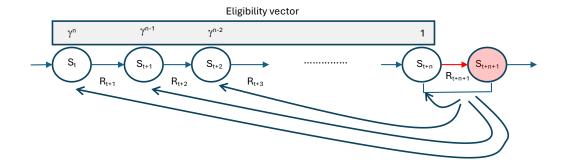
$$G_{t}^{n} = R_{t+1} + \gamma R_{t+2} \dots \gamma^{n} R_{t+n+1} + \gamma^{n+1} V^{n} (S_{t+n+1})$$

$$= \Sigma_{i=0}^{n} (\gamma^{i} R_{t+i+1} + \gamma^{i+1} V^{n} (S_{t+i+1}^{n}) - \gamma^{i} V^{n} (S_{t+i}))$$

$$= \Sigma_{i=0}^{n} \gamma^{i} (R_{t+i+1} + \gamma V^{n} (S_{t+i+1}) - V^{n} (S_{t+i}))$$

$$= \Sigma_{i=0}^{n} \gamma^{i} G_{t+i}^{0}$$

$$V^{n}(S_{t}) = (1-\alpha) V^{n}(S_{t}) + \alpha (G_{t}^{n})$$



## $TD(\lambda)$ : get the effect of TD(n) without having to pick an n

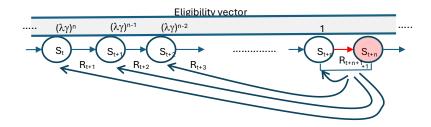
Intuition: rather than drop old states from sliding window, use exponential recencyweighted averaging trick

Sliding window for TD(n):

- Multiply Eligibility vector by  $\boldsymbol{\gamma}$
- Eligibility[ $S_{t+n+1}$ ]  $\leftarrow$  1
- Eligibility[ $S_t$ ]  $\leftarrow$  0

Sliding window for TD( $\lambda$ ): replacing traces version

- New hyper-parameter  $\lambda$
- Multiply Eligibility vector by  $\gamma\lambda$
- Eligibility[St+n+1] ← 1



# Implementation of $TD(\lambda)$

Algorithm 1:  $TD(\lambda)$  - estimating state-value function with eligibility traces.

```
import numpy as np
state_values = np.zeros(n_states) # initial guess = 0 value
eligibility = np.zeros(n_states)
lamb = 0.95 # the Lambda weighting factor
state = env.reset() # start the environment, get the initial state
# Run the algorithm for some episodes
for t in range(n_steps):
 # act according to policy
 action = policy(state)
 new_state, reward, done = env.step(action)
 # Update eligibilities
 eligibility *= lamb * gamma
  eligibility[state] += 1.0
 # get the td-error and update every state's value estimate
  # according to their eligibilities.
 td_error = reward + gamma * state_values[new_state] - state_values[state]
  state_values = state_values + alpha * td_error * eligibility
 if done:
    state = env.reset()
  else:
    state = new_state
```

#### Alistair Ries blog

https://amreis.github.io/ml/reinf-

learn/2017/11/02/reinforcement-learning-eligibility-traces.html

# Summary: concepts and keywords

Model-free methods

Deterministic vs. stochastic policies

Offline vs. online methods

Boot-strapping methods: TD(0), TD(n),  $TD(\lambda)$ , SARSA, Q-learning

Non-boot-strapping methods: Monte Carlo

Bias-variance tradeoff between TD and Monte Carlo

Eligibility traces: implementation

