

RLearning:

Short guides to reinforcement learning

Unit 3-3: Temporal Difference Learning

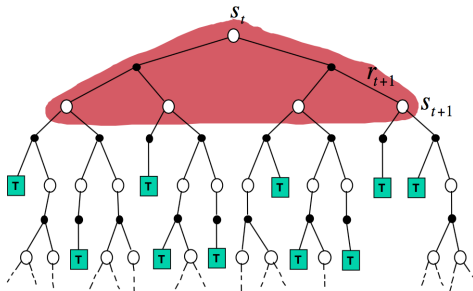
Davud Rostam-Afschar (Uni Mannheim)

How can we learn by sampling
from each step?

RL Algorithms

Dynamic Programming Backup

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

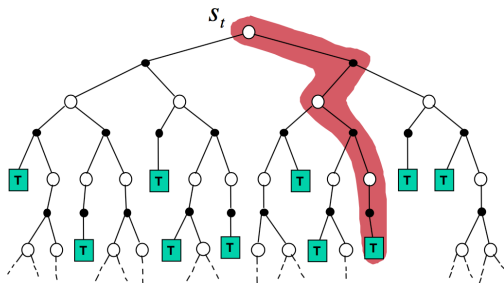


Source: David Silver

RL Algorithms

Monte Carlo Backup

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

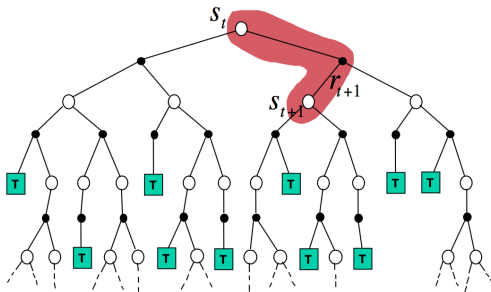


Source: David Silver

RL Algorithms

Temporal Difference Backup

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



Source: David Silver

Model Free Evaluation

- ▶ Given a policy π estimate $V^\pi(s)$ without any transition or reward model
- ▶ **Temporal difference** (TD) evaluation

$$V^\pi(s) = E[r \mid s, \pi(s)] + \gamma \sum_{s'} \mathbb{P}(s' \mid s, \pi(s)) V^\pi(s')$$
$$\approx r + \gamma V^\pi(s') \quad (\text{one draw approximation})$$

Toy Maze Example

3	r	r	r	+1
2	u		u	-1
1	u	l	l	l
	1	2	3	4

Start state: (1,1)

Terminal states: (4,2), (4,3)

No discount: $\gamma=1$

Reward is -0.04 for non-terminal states

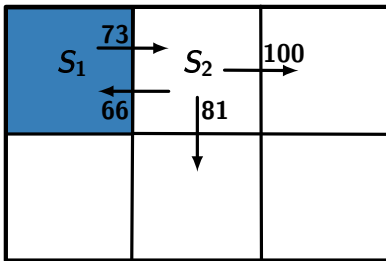
Four actions:

- ▶ up (**u**),
- ▶ left (**l**),
- ▶ right (**r**),
- ▶ down (**d**)

Do not know the transition probabilities

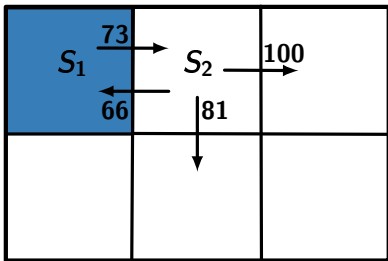
What is the value $V(s)$ of being in state s

Temporal Difference



$\gamma = 0.9, \alpha = 0.5, r = 0$ for non-terminal states

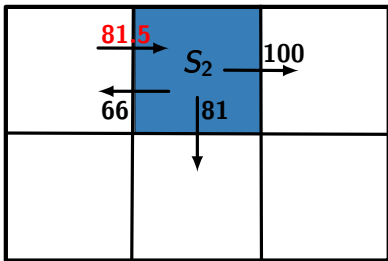
Temporal Difference



$\gamma = 0.9, \alpha = 0.5, r = 0$ for non-terminal states

$$\begin{aligned} Q(s_1, right) &= Q(s_1, right) + \alpha \left(r + \gamma \max_{a'} Q(s_2, a') - Q(s_1, right) \right) \\ &= 73 + 0.5(0 + 0.9 \max\{66, 81, 100\} - 73) \\ &= 73 + 0.5(17) \\ &= 81.5 \end{aligned}$$

Temporal Difference



$\gamma = 0.9, \alpha = 0.5, r = 0$ for non-terminal states

$$\begin{aligned} Q(s_1, right) &= Q(s_1, right) + \alpha \left(r + \gamma \max_{a'} Q(s_2, a') - Q(s_1, right) \right) \\ &= 73 + 0.5(0 + 0.9 \max\{66, 81, 100\} - 73) \\ &= 73 + 0.5(17) \\ &= 81.5 \end{aligned}$$

Temporal Difference Evaluation

Temporal Difference Evaluation

- ▶ Approximate value function: $V_n^\pi(s) \approx r + \gamma V^\pi(s')$
- ▶ **Incremental update** of sample (π, s', s)

$$V_n^\pi(s) \leftarrow V_{n-1}^\pi(s) + \alpha_n (r + \gamma V_{n-1}^\pi(s') - V_{n-1}^\pi(s))$$

Exploration vs Exploitation

Stochastic approximation (Robbins-Monro algorithm)

- ▶ **Theorem:** If α_n is appropriately decreased with number of times a state is visited then $V_n^\pi(s)$ converges to correct value
- ▶ **Sufficient conditions** for α_n :

$$\sum_n \alpha_n \rightarrow \infty \quad (1)$$

$$\sum_n \alpha_n^2 < \infty \quad (2)$$

- ▶ Often $\alpha_n(s) = 1/n(s)$,
where $n(s) = \#$ of times s is visited

Temporal Difference (TD) Evaluation

TD evaluation (π, V^π)

Repeat

Execute $\pi(s)$

Observe s' and r

Update counts: $n(s) \leftarrow n(s) + 1$

Learning rate: $\alpha \leftarrow 1/n(s)$

Update value: $V^\pi(s) \leftarrow V^\pi(s) + \alpha (r + \gamma V^\pi(s') - V^\pi(s))$

$s \leftarrow s'$

Until convergence of V^π

Return V^π

Temporal Difference Control

Temporal Difference Control

- Approximate Q-function:

$$\begin{aligned} Q^*(s, a) &= E[r \mid s, a] + \gamma \sum_{s'} \mathbb{P}(s' \mid s, a) \max_{a'} Q^*(s', a') \\ &\approx r + \gamma \max_{a'} Q^*(s', a') \end{aligned}$$

- **Incremental update**

$$Q_n^*(s, a) \leftarrow Q_{n-1}^*(s, a) + \alpha_n (r + \gamma \max_{a'} Q_{n-1}^*(s', a') - Q_{n-1}^*(s, a))$$

Comparison

- ▶ **Monte Carlo evaluation:**
 - ▶ Unbiased estimate
 - ▶ High variance
 - ▶ Needs many trajectories
- ▶ **Temporal difference evaluation:**
 - ▶ Biased estimate
 - ▶ Lower variance
 - ▶ Needs less trajectories

References I

- DENERO, J., D. KLEIN, B. MILLER, N. HAY, AND P. ABBEEL (2013): "The Pacman AI Projects," <http://inst.eecs.berkeley.edu/~cs188/pacman/>, Developed at UC Berkeley. Core by John DeNero and Dan Klein; student autograding by Brad Miller, Nick Hay, and Pieter Abbeel.
- POUPART, P. (2025): "Pascal Poupart's Homepage," <https://cs.uwaterloo.ca/~ppoupart/>, Accessed: 2025-05-24.
- RUSSELL, S. J., AND P. NORVIG (2016): *Artificial intelligence: a modern approach*. Pearson.
- SIGAUD, O., AND O. BUFFET (2013): *Markov decision processes in artificial intelligence*. John Wiley & Sons.
- SUTTON, R. S., AND A. G. BARTO (2018): "Reinforcement learning: An introduction," *A Bradford Book*, Available at <http://incompleteideas.net/book/the-book-2nd.html>.
- SZEPESVÁRI, C. (2022): *Algorithms for reinforcement learning*. Springer nature, Available at <https://sites.ualberta.ca/~szepesva/RLBook.html>.

Takeaways

How Does TD Learning Update Value Estimates Step-by-Step?

- ▶ No need to know transition probabilities or reward function
→ Model free
- ▶ Combine immediate reward with estimated value of next state
→ Biased value estimation from bootstrapped samples
- ▶ Revises estimates after each observed step
→ Needs few trajectories
- ▶ Lower variance than Monte Carlo at the cost of some bias
→ Bias-variance tradeoff
- ▶ Often used in algorithms such as Q-learning and SARSA