

Model-Free Q-learning with MC and TD ²

DOROZHKO Anton

Novosibirsk State University

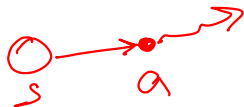
May 20, 2020

²David Silver's Lecture 4

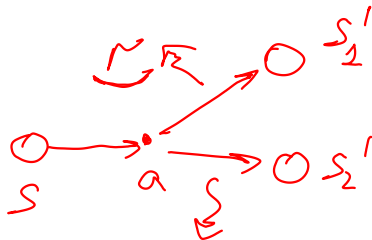
Outline

$$Q(s, a) = \mathbb{E}[G]$$

1 Markov decision processes



2 PI vs VI



3 MC and TD

Markov Decision Process

Definition

A Markov Decision Process is a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$

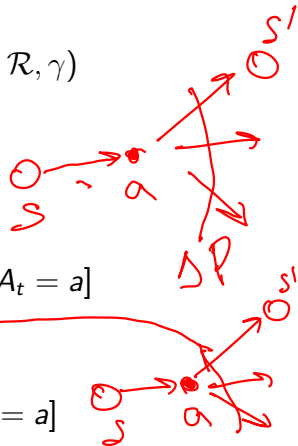
- \mathcal{S} is a (finite) set of states
- \mathcal{A} is a finite set of actions
- \mathcal{P} is a transition probability matrix

$$\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$$

- \mathcal{R} is a reward function :

$$\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$$

- γ is a discount factor, $\gamma \in [0, 1]$ *Horizon*



Recall Policy Iteration

Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Loop:

$\Delta \leftarrow 0$

Loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

fixed $\pi \rightarrow V^\pi$

3. Policy Improvement

policy-stable \leftarrow true

For each $s \in \mathcal{S}$:

old-action $\leftarrow \pi(s)$

$\pi(s) \leftarrow \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

If *old-action* $\neq \pi(s)$, then *policy-stable* \leftarrow false

If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

$V^\pi \rightarrow \pi_{i+1}$

Recall : Value Iteration

Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

| $\Delta \leftarrow 0$

| Loop for each $s \in \mathcal{S}$:

| $v \leftarrow V(s)$

| $V(s) \leftarrow \max_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$

| $\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that

$\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$

Comparison ⁸

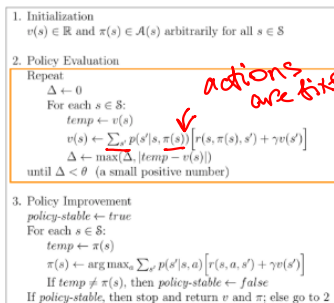


Figure 4.3: Policy iteration (using iterative policy evaluation) for v_* . This algorithm has a subtle bug, in that it may never terminate if the policy continually switches between two or more policies that are equally good. The bug can be fixed by adding additional flags, but it makes the pseudocode so ugly that it is not worth it. :-)

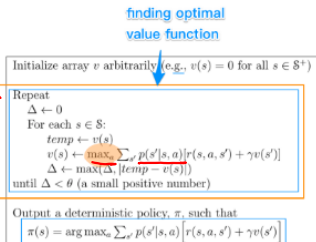


Figure 4.5: Value iteration.

one policy update (extract policy from the optimal value function)

$$\underline{p(s'|s, a)}$$

Monte-Carlo RL

- MC learns directly from episodes of experience
- Model-free: no knowledge of

$$\mathcal{P}_{ss'}^a \quad \text{or} \quad \mathcal{R}$$

- complete episodes
- idea : value = mean return
- Caveat: only episodic MDPS **episodes must terminate**

MC Evaluation

$$V_{\pi}(s)$$

$$Q_{\pi}(s, a)$$

- Goal: learn V_{π} from episodes under policy π

$$S_1, A_1, R_2, \dots, S_k \sim \pi$$

- Total discounted reward

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

- Value function

$$V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s]$$

$$\frac{1}{n} \sum G_t$$

MC Evaluation

- Goal: learn V_π from episodes under policy π

$$S_1, A_1, R_2, \dots, S_k \sim \pi$$

- Total discounted reward

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

- Value function

$$V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s]$$

- Monte-Carlo policy evaluation uses empirical mean return instead of expected return

First-Visit Monte-Carlo Policy Evaluation

$$s, a, r, s_1, a_1, r_1 = \dots \sim \pi$$

$$\begin{matrix} (s \\ s \end{matrix}$$

$$G$$

- To evaluate state s
- The **first** time-step t that state s is visited in each episode

- $\underline{N}(s) = N(s) + 1$ # of 1st appearance of s in each game

- $\underline{S}(s) = S(s) + \underline{G}_t$ - \sum over many games (episode)

- $\underline{V}(s) = S(s)/N(s)$

- By law of large numbers $V(s) \rightarrow V^\pi(s)$ as $N(s) \rightarrow \infty$

$$E[G_t | s]$$

Incremental Mean

$$V^\pi \rightarrow \pi$$

$$\mu_k = \frac{1}{k} \sum_{j=1}^k x_j$$

played k games
update policy

$$= \frac{1}{k} \left(x_k + \sum_{j=1}^{k-1} x_j \right)$$

$$= \frac{1}{k} (x_k + (k-1)\mu_{k-1})$$

$$= \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1})$$

from $k-1$
games

diff from
last game

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

$$N(S_t) = N(S_t) + 1$$

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

incremental
mean
 $\hat{V}(S_t)$

- In non-stationary problems: running mean

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

cum reward for episode
for finished
game

Temporal-Difference Learning

- TD learns from episodes
- model-free $P_{ss'}^a$ R
- **incomplete episodes**, bootstrapping
- Update a guess towards a guess
idea

don't want to wait until the end of episode

diff \leftrightarrow MC complete episodes

MC and TD

$$S_0, a_0, r_0, S_1, a_1, r_1 \leftarrow t$$

- Goal: learn V_π online from experience under policy π
- Incremental every-visit MC
 - $V(S_t) = V(S_t) + \alpha(G_t - \hat{V}(S_t))$

predictions (pointing to $\hat{V}(S_t)$)
- Simplest TD : TD(0)
 - Update $V(S_t)$ toward estimated return $R_{t+1} + \gamma V(S_{t+1})$

$$V(S_t) = V(S_t) + \alpha(\underbrace{R_{t+1} + \gamma V(S_{t+1})}_{G_t}) - V(S_t)$$

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

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

Driving Home Example

State	Elapsed Time (minutes)	TD Predicted Time to Go	MC Predicted Total Time
leaving office	0	<u>30</u>	<u>30</u>
reach car, raining	5	<u>35</u> <i>ETA</i>	40
<u>exit highway</u>	20	15	35
<u>behind truck</u>	30	10	40
home street	40	<u>3</u>	43
arrive home	43	<u>0</u>	<u>43</u>

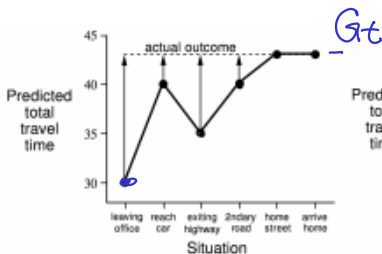
*Get**Whole traj*

MC vs TD

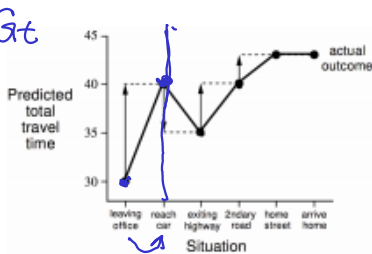
$$V(S_{t+1}) = V(S_t) + \alpha (G_t - V(S_t))$$

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

Changes recommended by Monte Carlo methods ($\alpha=1$)



Changes recommended by TD methods ($\alpha=1$)



Advantages and Disadvantages

- TD can learn before knowing the **final outcome**
 - TD learn after each step
 - MC must end the episode
- TD can learn without the final outcome
 - TD - incomplete sequences
 - TD - continuing envs
 - MC - complete sequences
 - MC - only episodic envs

Handwritten diagram illustrating the TD learning process:

A horizontal timeline axis is shown with several vertical tick marks representing time steps. Above the axis, there are two upward-pointing arrows, one at the second tick mark and one at the fifth tick mark, with the label ± 1 to the right of the second arrow. Below the axis, the expression $R_{t+1} - V(S_{t+1}) - V(S_t)$ is written. A bracket underneath the terms $-V(S_{t+1}) - V(S_t)$ is labeled with the Greek letter δ . The label $G_t = 0$ is written in the upper right corner.

Advantages and Disadvantages (2)

MC has high variance, zero bias

- Good convergence
- Not sensitive to initial value
- Simple

TD has low variance, some bias

- More efficient than MC
- TD(0) converges to V_π
- (not always for function approximation)
- Sensitive to initial value

$$\sum R \leftarrow r.V$$

$$G_t - V(S_t)$$

$$TD(0) \quad R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

$$R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2}) - V(S_t)$$

TD(1)

$$R + \gamma V(S_{t+1}) - V(S_t)$$

Advantages and Disadvantages (3)

$$p(s_{t+1} | s_t) = p(s_{t+1} | s_t, \mathbf{x}_t) \quad \text{X}$$

TD exploits Markov property

- Usually more efficient in Markov envs

MC does not exploit Markov property

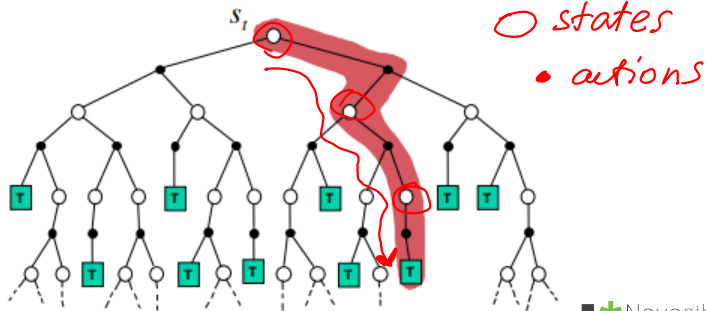
- Usually more effective in non-Markov envs

play the whole game
↓
count

lose ordering

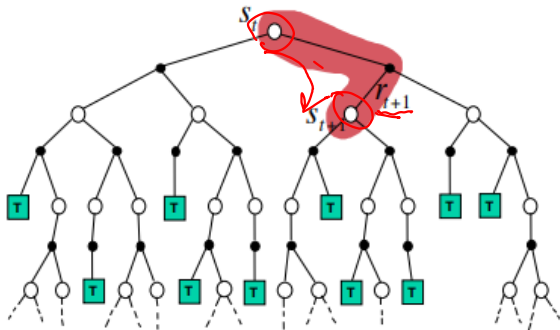
MC Backup

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



TD Backup

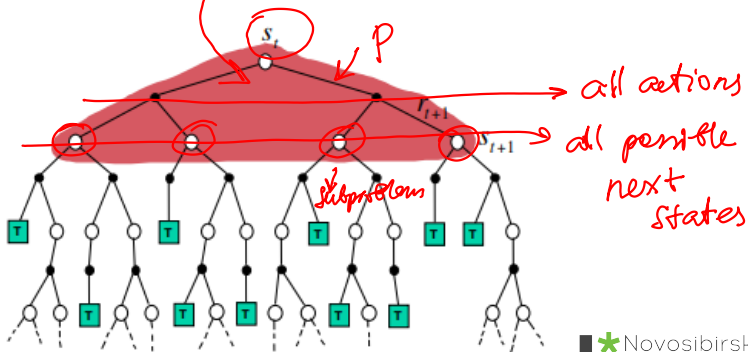
$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



DP Backup

Bellman Expectation eq.

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



Bootstrapping and Sampling

Bootstrapping update involves an estimate

- MC
- DP
- TD

Sampling update sample an expectation

- MC
- DP
- TD

Bootstrapping and Sampling

Bootstrapping update involves an estimate

- MC ✗
- DP ✓
- TD ✓

Sampling update sample an expectation

- MC
- DP
- TD

Bootstrapping and Sampling

Bootstrapping update involves an estimate

- MC ✗
- DP ✓
- TD ✓

Sampling update sample an expectation

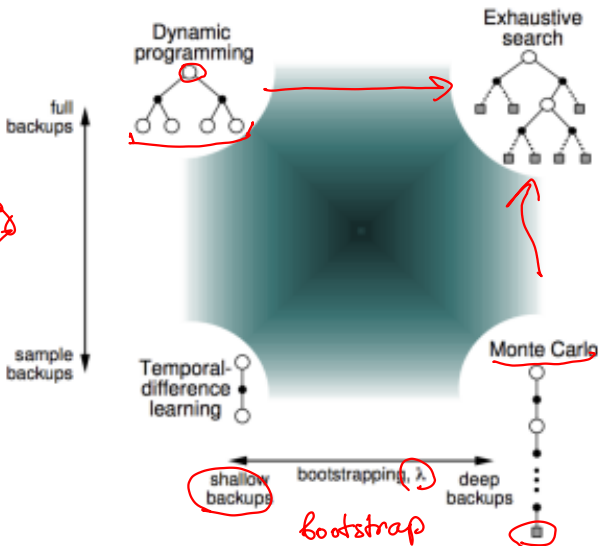
- MC ✓
- DP ✗
- TD ✓



Unified View of RL

Model Free

~~DP~~ ~~ES~~



brute force
 $E \sim \frac{1}{n} \Sigma$

"sample"
 played
 so many
 games

Value and Policy Iteration Lab



<https://bit.ly/2JVv6rc>