

Bootstrapping:

It is a technique used in RL to estimate the value of a state or state-action pair by using the agent's current estimate of the value func<sup>n</sup>.

→ This is done by updating the value func<sup>n</sup> using the TD error, which is the difference between the agent's current estimate value func<sup>n</sup> and the observed reward.

Why useful:

1. Efficiency
2. Scalability.
3. Better convergence
4. Learn from incomplete information.
5. Improve exploration - exploitation trade-off.
6. Reduce the impact of noise & uncertainty.

How it works:

- estimating the value func<sup>n</sup>  $V(s)$  which is the expected return from state  $s$ .
- Bootstrapping means that we use the current estimate of  $V(s')$    
 value of next state   
 to update  $V(s)$
- This is typically done using TD learning, where updates happen in real-time as the agent moves through the environment.

TD learning:

- Includes TD(0), SARSA & Q-learning.
- general update rule in TDL for a state " $s$ " at time " $t$ " is

$$V(s) \leftarrow V(s) + \alpha \cdot (r + \gamma V(s') - V(s))$$

$\gamma$  = discount factor

$\alpha$  = Alpha is learning rate

$V(s')$  = estimated next state

$V(s)$  = <sup>estimated</sup> current state



Algorithm uses in BS:

1. Q-learning
2. SARSA
3. Deep Q-Networks (DQN)

} uses TD bootstrapping to update the action-value function.

Adv:

- Faster updates & more efficient learning.
- Provides more frequent learning signals.

Disadv:

- BS introduces "bias" because it relies on estimates rather than actual outcomes.
- This can make learning process less stable.

TD(0) Algorithm:

- TD is a technique in RL that combines both Monte Carlo methods and Dynamic programming.
- It is a model-free algorithm which does not require any information or knowledge of the environment.
- In this method Agent will update its value estimates immediately after taking an action, based on the current state/reward & estimated value of the next state.
- In TD(0) "0" represents single-step update. means the value of a state is updated based on the next immediate state's value, without looking further states.
- TD(0) combines bootstrapping feature of DP and the sampling-based approach of MC.
- It focuses only on-policy evaluation.

Objective:

Main goal of TD(0) is to estimate the value function  $v(s)$  for given policy  $\pi$  where

$$v(s) = E_{\pi}[G_t | s_t = s]$$

$G_t$  = return.



### TD(0) update rule:

In this the value will be updated after each step by combining the current state " $s$ " & next state " $s'$ ".

→ state value function

$$V(s) \leftarrow V(s) + \alpha (r + \gamma V(s') - V(s))$$

$V(s)$  = estimate value of current state

$V(s')$  = " " " next "

$\gamma$  = discount factor

$\alpha$  = Alpha is a learning rate

$r$  = immediate reward

### TD error:

It represents the difference b/w current state and updated target, which has reward plus for next state.

→ TD error for a state  $s$  after observing reward  $r$  & transitioning to state  $s'$  is

$$\delta = r + \gamma V(s') - V(s)$$

### updated rule:

$$V(s) \leftarrow V(s) + \alpha \cdot \delta$$

### Algorithm steps:

1. Initialize.

$V(s)$  for all states  $s$ .

2. Repeat.

• for each episode

i. Take action — based on policy  $\pi$  & observe next state  $s'$  & reward  $r$

ii. compute TD error —  $\delta = r + \gamma V(s') - V(s)$

iii. update the value :  $V(s) \leftarrow V(s) + \alpha \cdot \delta$



iv. Move to next state  $s'$

5. Repeat until convergence.  
(or) set of episodes are completed.

Characteristics of TD(0):

- on-policy
- Incremental updates
- Bootstrapping

Benefits:

- Efficiency
- Suitability for continuous task
- Less memory requirement

Limit:

- Dependency on Discount factor " $\gamma$ "
- Bias in estimation.

Apply:

1. Robotics / Robot navigation
2. game playing

Convergence of Monte Carlo & batch TD(0) Algorithms:

MC & TD are fundamental methods in RL for estimating value function.  
Both algorithms are evaluate the value function  $V^\pi(s)$  for a given policy  $\pi$ .

Monte Carlo:

Imagine playing a game and making moves randomly.

→ After finishing the game, you look at the outcome and assign values to each move based on whether it led to a win or loss.

→ you update your strategies for each move based on the overall result of the entire game.

Convergence of MC:

update:

$$V(s) \leftarrow V(s) + \alpha \cdot (G_t - V(s))$$

$G_t$  = observed reward from state  $s$ .

Key convergence:

1. consistency
2. weakness in Bootstrapping
3. unbiased estimate



batch TD(0):

In this Algorithm Instead of waiting until the end of the game you update your strategy after each move.

→ you estimate the value of each move by comparing it to the value of the next move.

→ This helps you learn and adjust your strategy as you go through the game.

update rule:

$$\text{TD error is } \hat{V}(s) \leftarrow V(s) + \alpha \delta_t V(s)$$

$$\delta_t = r + \gamma V(s') - V(s)$$

Batch mode:

TD(0) updates are applied after collecting all state transitions in a batch.

• Batch TD(0) minimizes the MSE mean square error b/w the predicted value  $V(s)$  & the BS target

$$\text{MSE} = \sum_s [\gamma + \gamma V(s') - V(s)]^2$$

convergence of batch TD(0):

1. fixed-point convergence
2. linear equ<sup>n</sup> framework
3. conditions for convergence.
4. faster updates.

MC	batch TD(0)
<ol style="list-style-type: none"><li>1. It is better for episodic tasks where episodes are short &amp; complete return information is readily available.</li><li>2. It is computationally expensive &amp; slow in practice</li></ol>	<ol style="list-style-type: none"><li>1. It is better for continuous tasks or long tasks.</li><li>2. It is faster, but with a slight bias introduced by bootstrapping.</li></ol>



## Model-Free Algorithm / control.

MFC is a type of RL algorithm that does not require an explicit model of the environment to "create" a policy that maximizes the sum of future rewards.

### features:

1. No model of the environment.
2. Focus on optimal policy.
3. Exploration vs Exploitation
4. Common Algorithms.

- i. value-based methods
- ii. policy-based "
- iii. Actor-critic "

### i. value-based methods:

learn action-value func<sup>n</sup>  $Q(s,a)$  that estimates the expected return  
→ it has two approaches

#### a. Q-learning:

- off-policy methods that learn  $Q(s,a)$  using the Bellman optimality eqn<sup>n</sup>.

$$Q(s,a) \leftarrow Q(s,a) + \alpha [r + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$

updates are based on best possible action.

#### b. SARSA

- state-action-reward-state-action
- on-policy method

$$Q(s,a) \leftarrow Q(s,a) + \alpha [r + \gamma Q(s',a') - Q(s,a)]$$



ii. policy-based methods:

Directly learn the policy  $\pi(a|s) \setminus p(a|s) \pi(a|s)$  without estimating a value function.

Reinforce Algorithm:

• A MC policy-gradient method:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log f(\pi_{\theta}(a_t|s_t)) G_t$$

$\downarrow$  policy parameter                       $\downarrow$  reward

iii. Action-critic method:

combine a value function with a policy (actor) for more stable updates.

Adv:

1. Simplicity.
2. Versatility.
3. Incremental learning

Chall:

1. Stability
2. Scalability

App:

1. Robotics
2. Health care
3. Autonomous vehicles
4. Games.

SARSA: (or) SARS'A'

State - Action - Reward -  $s'$  (Next State) - Action (for  $s'$ )

→ It is an on-policy RL Algorithm that learns the action-value function  $Q(s,a)$  based on Agent's current policy.

update rule:

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

$\alpha$  = learning rate

$\gamma$  = Discount factor

$Q(s',A')$  = Q value of next state & action.



→ SARSA updates the Q-value using action  $A'$  chosen according to the current policy making it on-policy.

Adv:

1. Balance exploration & exploitation.
2. safer in environment

Limitations:

1. May converge slower than off policy methods like Q-learning

Expected SARSA

It is a variant of SARSA that improves the stability of updates by using expected value of the Q-function at next state  $s'$  instead of relying on a single ~~action~~ action  $A'$ .

Update rule:

$$Q(s, A) \leftarrow Q(s, A) + \alpha \left[ R + \gamma \sum_{a'} \underbrace{\pi(a' | s')}_{\substack{\text{probability of taking} \\ \text{action } a' \text{ in state} \\ s' \text{ under} \\ \text{current policy}}} Q(s', a') - Q(s, A) \right]$$

Adv:

1. more stable.
2. improving converge in some scenarios.

Limit:

1. more expensive than SARSA

SARSA	Expected SARSA
1. on-policy	1. on-policy
2. uses a single sampled action $A'$	2. uses expected value over all actions $A'$
3. Higher variance	3. lower variance
4. less computationally expensive	4. more intensive
5. slower in some cases	5. faster & more stable



## Q-learning:

→ It is one of the most fundamental & widely used algorithms in RL.

→ It is a model-free

→ It is off-policy learning algorithm that can learn from other policies.

Q-learning is an off-policy LA that learns the optimal policy by directly estimating the optimal action-value func<sup>n</sup>  $Q^*(s, a)$  regardless of the current policy.

### Action-value func<sup>n</sup>:

$Q(s, a)$  represents reward (action), state  $s$ , taking action  $a$ .

$$Q(s, a) = E \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a \right]$$

$r_{t+k+1}$  = reward received at step  $t+k+1$

$\gamma$  = discount factor.



### goal:

~~optimal~~ To learn 'optimal action-value func<sup>n</sup>'  $Q^*(s, a)$  such that agent can determine the optimal policy.

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

### Q-learning Algorithm:

given state  $s$ , action  $a$ , Q-value is updated by following

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

$\alpha$  = learning rate

$\gamma$  = discount factor

$Q(s, a)$  = current estimate Action-value func<sup>n</sup>

$Q(s', a')$  = next state

$\max_{a'} Q(s', a')$  = maximum predicted Q-value for next state/action

$r$  = immediate reward.



Steps:

1. Initialization:

- Initialize  $Q(s, a)$   
state-action pair  $= (s, a)$   
 $\therefore Q(s, a) = 0$  Initially
- choose  $\alpha$  &  $\gamma$  parameters

2. Interaction with environment:

- observe current state  $s_t$
- select an action based on exploration strategy,

3. Execute the action:

perform action " $a_t$ ", observe reward " $r_t$ " then  
 $s_{t+1} = \{t+1\} s_{t+1}$

4. update the Q-value:

After all these update Q-value

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \cdot [r_t + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)]$$

5. Repeat:

until it converges

Adv:

1. model-free
2. off-policy
3. widely applicable

Limit:

1. Scalability
2. Exploration

Applica:

1. Game AI
2. Robotics
3. Finance
4. Healthcare.

Extensions of Q-Learn

1. Deep Q-learning - handle large datasets (or) continuous state spaces
2. Double Q-learning - Reduces overestimation.
3. Q-lambda