

Bootstrapping

- in RL refers to a technique when an agent updates its estimate of state or action values using current or previous estimates rather than relying solely on complete trajectories or final outcomes.
- This approach combines observed rewards with existing value predictions to refine its policy iteratively.

Mechanism

- In bootstrapping, the agent updates the value function $v(s)$ or $\varphi(s, a)$ using the Bellman

equations which express the relationship between the value of a state and the value of subsequent states.

for example

in Temporal Difference (TD) learning
the update rule for a state s_t is

$$v(s_t) \leftarrow v(s_t) + \alpha [r_{t+1} + \gamma v(s_{t+1}) - v(s_t)]$$

$\alpha \rightarrow$ learning rate.

$\gamma \rightarrow$ discount factor.

$r_{t+1} + \gamma v(s_{t+1})$ is bootstrapped estimate of the return.

TD(0) Algorithm

TD(0) (temporal difference learning) is a model-free, online RL algorithm that updates value estimates incrementally by combining observed rewards with current value predictions (bootstrapping).

It operates on a per time step basis without waiting for an entire Episode to complete.

TD(0) update Rule

The value of state $v(s_t)$ is updated using

$$v(s_t) \leftarrow v(s_t) + \alpha [r_{t+1} + \gamma v(s_{t+1}) - v(s_t)]$$

α - learning rate ($0 < \alpha \leq 1$)

γ - Discount factor ($0 \leq \gamma \leq 1$)

r_{t+1} - immediate reward after transitioning from s_t .

$v^*(s_{t+1})$ estimated value of the next state.

Algorithm

- initialize $v(s)$
- for each episode:
 - observe initial state s_t .
 - while s_t is not terminal.
 - take action a_t (using policy).
 - observe reward r_{t+1} and next state s_{t+1}
 - update $v(s_t)$ using the TD(α) rule
 - set $s_t \leftarrow s_{t+1}$

Ex:

current $V(s_2) = \sqrt{5} = 5$.

next state $s_2+1 = \sqrt{5+1} = 6$.

Observe reward

$$r_{t+1} = 2, \alpha = 0.1, \gamma = 0.9$$

Update

$$V(s_t) \leftarrow 5 + 0.1 [2 + 0.9 \times 6 - 5]$$

$$= 5 + 0.1 [2 + 5.4 - 5]$$

$$= 5 + 0.1 \times 2.4 = 5.24.$$

Advantages

- ① Efficiency
- ② low variance.
- ③ works in non-terminations environments.

Disadvantages

- ① Bias
- ② slower convergence.

Convergence of Monte Carlo and Batch TD(0) Algorithms

Both are used in RL to estimate value functions. Here, we focus on Monte Carlo and Batch TD(0) alg.

① Monte Carlo method

→ It learns value function by averaging the returns observed after visiting a state. They require complete episodes to update value estimates.

Convergence

mc methods converge to
the optimal value function
under conditions:

- sufficient exploration,
- infinite visits to each state-action pair.
- the policy evaluated is fixed

Advantages

- unbiased estimates because they use actual return
- simple to implement

Disadv

- requires complete episode before updates.
- high variance

③ Batch TD(0) method

TD(0) is a bootstrapping method that updates value estimate based on the observed reward and the estimated value of the next state.

Batch TD(0) applies TD(0) update to a fixed dataset (batch) of experiences.

Convergence

Batch TD(0) converges to the value function that minimizes the mean squared error on the training data.

→ The batch of data is fixed and finite.

→ The markov decision process (MDP) is finite and deterministic.

→ the step-size parameter
(learning rate) satisfies the
Robbins - monk conditions.

Advantages

- can update value estimate
after each step (online
learning).
- lower variance compared
to mc method.

Disadvan

- Biased estimate due to
bootstrapping
- Requires careful tuning
of the learning rate.

Key diff

Batch TD(0)

Monte Carlo

converges to set,
optimal value function

value function
minimizations
MSE on batch.

Blau variance
unbiased, high
variance

Biased, lower
variance.

Data Requirement
Requires complete
episodes

can use incomplete
episodes

Speed of Learning
Slower due to high
variance

Faster due to
bootstrapping

① Monte Carlo method

you're playing a game where you move through rooms and collect rewards. The goal is to estimate the value of each room (state) based on the total reward you collect.

States: Room A, Room B, Room C.

Actions: move left, move right

Rewards

→ moving from Room A to Room B gives +1 reward.

→ moving from Room A to Room C gives +1 reward.

→ Room B and Room C are terminal states (game end).

Episode - 1

Path:- $A \rightarrow C$ (reward = 1)

Total return from Room A

$$G(A) = 1$$

Episode - 2

Path:- $A \rightarrow B$ (reward = 0)

Total return from Room A:

$$Q(A) = 0$$

Episode - 3

Path:- $A \rightarrow C$ (reward = 1)

Total return for Room A:

$$G(A) = 1$$

Value Estimation for Room A

$$V(A) = \frac{1+0+1}{3} = \frac{2}{3}$$

TD(0) method

Batch of Experience

(A, right, 1, C)

(A, left, 0, B)

(A, right, 1, C)

Step 1: initialize $v(A) = 0, v(B) = 0, v(C) = 0$

Step 2: for each experience, compute

the TD target:

$$TD_{target} = r + \gamma v(s')$$

Assume discount factor $\gamma = 1$)

For (A, right, 1, C)

$$TD_{target} = 1 + 1 \cdot v(C) = 1 + 0 = 1$$

for (A, ~~right~~ left, 0, B):

$$TD_{target} = 0 + 1 \cdot v(B) = 0 + 0 = 0$$

For (A, right, 1, C):

$$TD_{target} = 1 + 1 \cdot v(C) = 1 + 0 = 1$$

update $v(A)$ at the arrival
of 2D dataset

$$v(A) = \frac{1+0+1}{3} = \frac{2}{3}.$$

Model-free control in RL

* Model-free control in RL refers to a class of algorithms where the agent learns to make optimal decisions without explicitly constructing a model of the environment. Instead, the agent directly interacts with the environment, learns from the experience (state, action, reward, next state), and improves its policy over time.

Common model-free control

Algorithms

① Monte Carlo Methods

- learn value functions or policies by averaging returns from complete episodes.
- suitable for episodic tasks.
- Example:- Monte carlo control, where the agent updates the policy after each episode.

② Temporal Difference (TD) methods

- learn value functions or policies by bootstrapping from current estimate.
- suitable for both episodic and continuous tasks.

→ Example

→ SARSA

→ φ -learning.

③ Policy Gradient Method

→ Directly optimizes the policy by maximizing the expected return using gradient descent

$$\hat{S}_t \leftarrow$$

→ REINFORCE Algorithm

→ Actor-Critic Method.

Advantages

- ① Simplicity.
- ② Applicability
- ③ Scalability.

Disadvantages

- ① Simple inefficiency
- ② Exploration challenges
- ③ Convergence issues

Q-Learning

- * Q-learning is a model free, off-policy reinforcement learning algorithm that learns the optimal action-value function $q(s, a)$, representing the expected utility of taking action a in state s .
- * goal is to maximize cumulative reward by iteratively updating a q -table, which guides the agent's decisions.

Key components

Q-table: Matrix storing $q(s, a)$ values for state-action pairs

Bellman Equation :- Basis for q -value update.

$$q(s, a) \leftarrow q(s, a) + \alpha [R(s, a) + \gamma \max_{a'} q(s', a') - q(s, a)]^{\alpha}$$

α : learning rate ($0 \leq \alpha \leq 1$)

γ : discount factor

$r(s,a) \rightarrow$ immediate reward
after taking action a
in state s

$\max_a q(s', a') \rightarrow$ maximum
feature reward
estimate from next state s' .

Algorithm

- ① initialize q -table \rightarrow Typically with zeros or small random values.
- ② select action \rightarrow use an exploration strategy to choose action a in state s ,
- ③ execute action \rightarrow observe reward r and next state s' ,

④ Update φ -values: Apply the Bellman equation to adjust $\varphi(s,a)$

⑤ Repeat: Transition to s' and iterate until convergence or terminal state.

Ex:- Simple Grid world

Environment

States: s_0 (start), s_1 (goal)

Actions: Right (R), Down (D)

(but Down D is invalid in s_0)

Rewards: $s_0 \rightarrow s_1$ via Right + 1
(goal reached)

→ invalid operation

(Down in s_0): 0 (no movement)

Terminal state: s_1 (episode end)

Step-1: initializing Q-table

state	right (R)	down (D)
s ₀	0	0
s ₁	0	0

Step 2: Agent interaction (One Epoch)

Parameters

- learning rate (α) = 0.1
- discount factor (γ) = 0.9
- Exploration (ϵ) = 0.1 (10% chance of random action),

Episode Flow

- ① Start at s₀
- ② choose Action
 - with 10% chance (ϵ), explore: Down (D) invalid
 - with 90% chance, exploit: right (R) (optimal)

Assume the agent explores and selects Down (D) first.

(2) Observe Reward

- Reward = 0 (no movement)
- next state = s_0 (stays in place)

(3) Update ϕ -value for $s_0 \rightarrow D$

$$\phi(s_0, D) = \phi(s_0, D) + \alpha [R + \gamma \max_{a'} \phi(s_0, a') - \phi(s_0, D)]$$

$$\phi(s_0, D) = 0 + 0.1 [0 + 0.9 \times \max(0, 0) - 0]$$

$$\phi(s_0, D) = 0$$

No change since s_0 has no future reward yet.

(4) next step

→ ~~Agent~~ tries again in s_0

→ Now, Exploit Right (R)
(Greedy step)

6. Observe Result

→ Reward = +1 (goal reached)

→ Next state = s_1 , (terminal).

④ Update q-value for $s_0 \rightarrow R$:

$$\varphi(s_0, R) = \varphi(s_0, R) + \alpha [R + r \max_a \varphi(s_0, a) - \varphi(s_0, D)]$$

$$\varphi(s_0, R) = 0 + 0.1 [1 + 0.9 \times \max_a (\varphi(s_1, a)) - 0]$$

$$\varphi(s_0, R) = 0 + 0.1 \times 1 = 0.1$$

Update φ -table

state	Right(R)	Down(D)
s_0	0.1	0
s_1	0	0

Final ϕ -table (After convergence)

State	Right (R)	Down (D)
s_0	1.0	0
s_1	0	0

The agent learns that right from s_0 is optimal to reach the goal.

$$\text{Episode 1} \rightarrow 0.1$$

$$\text{Episode 2} \rightarrow 0.19$$

...

$$\text{Episode } n \rightarrow 1.0$$

Exploration vs Exploitation

① Exploration

* trying new actions to discover their reward

* Helps the agent learn about the environment.

② Exploitation

- uses known information to choose the best action.
- maximizes immediate reward on current knowledge.

Trade-off

There should be no 'too much Exploration' and 'too much Exploitation'

Exploration

if Too much Exploration → Agent wastes time on bad actions.

if too much Exploitation →

Agent may miss better actions.

Exploration → trying a new path to find shortcut

Exploitation → trying the known fastest path.

SARSA

→ SARSA (State - Action - Reward -

Action) is an on-policy π^*

learning algorithm, it learns the φ -values $\varphi(s, a)$ for a policy while following the same policy, incorporating exploration into updates. The name reflects the sequence of steps:

(s, a, r, s', a')

key components

→ on-policy learning update φ value based on actions taken by the current policy

→ update rule

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

$\alpha \rightarrow$ learning rate

$\gamma \rightarrow$ discount factor

$\varphi(s', a') \rightarrow$ φ -value of the next action selected by the policy

3. Algorithm steps

1. Initialize $\phi(s, a)$ arbitrarily
(e.g. zero)
2. for each episode
 - a. Start ~~won~~ in state s .
 - b. choose action a using the current policy (e.g., ϵ -greedy)
 - c. Repeat until terminal state:
 - take action a , observe reward R , and next state s' .
 - choose next action a' using the current policy.
 - update $\phi(s, a)$ using the SARSA rule.
 - $s \leftarrow s'$, $a \leftarrow a'$

SARSA vs ϕ -learning

SARSA

on-policy
(learn the policy
being used)

uses $\phi(s', a')$
(actual next action)

Exploration

Accounts for
exploration in updates

~~Risk sensitivity~~
Safer in risky
environments

ϕ -learning

off-policy
(learns optimal
policy).

uses Max^{a'}
 $\phi(s', a')$.

ignores
exploration
in update.

more
aggressive

Ex:
Same for ϕ -learning Exa
Same formula just use

$$\phi(s_0, D) = \phi(s_0, D) + \alpha [r + \gamma \phi(s_0, R) - \phi(s_0, D)]$$

Expected SARSA

→ it is on-policy TD learning algorithm. It extends SARSA by using the expected value of the Q-value for the next state, weighted by the current policy action, probabilities. This reduces variance compared to SARSA, which relies on a single sampled action.

② Key components

- On-Policy Learning
- Expected value
- Update Rule

$$\varphi(s, a) \leftarrow \varphi(s, a) + \alpha [R_t + \gamma \sum_{a'} \pi(a'|s') \varphi(s', a') - \varphi(s, a)]$$

$\pi(a'|s')$ → probability of choosing action a' in state s' under the current policy.

Algorithm

- ① Initialize $\varphi(s, a) \rightarrow \text{zero}.$
- ② for each episode.
 - a) start in state s .
 - b) choose action a using the current policy
 - c) repeat until terminal state:
 - take action a , observe reward R_t and next state s' .
 - compute the expected φ -values for s' .

$$E[\phi(s', \cdot)] = \sum_a \pi(a'|s') \phi(s', a')$$

→ Update $\phi(s, a)$:

$$\phi(s, a) \leftarrow \phi(s, a) + \alpha [R + \gamma E[\phi(s', \cdot)] - \phi(s, a)].$$

→ choose next action a' using the policy.

$$\rightarrow s \leftarrow s', a \leftarrow a'.$$

Ex: Same Example a)
the φ-learner but
use formulae.

Agent is so choose Down (D)
Reward α , Next state so
Expected φ-value for s .

$$E[\varphi(s_0, \cdot)] = (0.9 \times 0) + (0.1 \times 0) = 0$$

update $\varphi(s_0, D)$:

$$\begin{aligned}\varphi(s_0, D) &= 0 + 0.1[0 + 0.9 \times 0 - 0] \\ &= 0\end{aligned}$$

$$\left. \begin{array}{l} \alpha = 0.1 \\ r = 0.9 \\ \epsilon = 0.1 \end{array} \right\}$$

so for right (R)
Reward +1
Expected φ -value
for $s_1 = 0$ (Terminal
state).

Expected SARSA	SARSA	φ -learning
Expectation over policy's actions	sampled next action.	max φ -value of next state
lower (user) expectation	Higher (averagingly sampled)	Highest (user man operator).
on-policy	on policy	off-policy

Short

① Episode

- A sequence of interactions from an initial state to a terminal state

Ex.: A robot navigation from start to goal, ends when it succeeds or fails.

② Policy

- The agent's strategy for choosing actions in a state

Type

- Deterministic: $\pi(s) = a$ (fixed action for a state)
- Stochastic: $\pi(a|s) = \text{probability of action } a \text{ in state } s$

③ Convergence

→ when the algorithm stabilizes
(φ -values or policy stop changing significantly).

Eg: φ -learning converges to the optimal φ -table after many episodes.

④ on-policy vs off-policy

on-policy

learns the policy being followed
(eg, SARS α)

updates w.e action from the current policy

Eg: SARS α

off-policy

learns the optimal policy independently of actions taken.
(eg: φ -learning)

update w.e hypothetical actions (max φ value).

Eg: φ learning.