

UCS2H24 PRINCIPLES OF REINFORCEMENT LEARNING

Unit 3

Integration of Tabular Methods

n-step Bootstrapping: TD Prediction - SARSA - OFF policy learning;
Planning and Learning with Tabular Methods

* Temporal Difference (TD) Learning

- a method in RL aimed at predicting future outcomes or optimizing decisions in uncertain environments.
- combines aspects of:
 - (i) Monte Carlo methods - estimate outcomes based on complete episodes
 - (ii) Dynamic Programming - uses known ~~values~~^{models} of the environment.
- The primary goal of TD learning is to estimate the value of a state/action in a given environment. It involves updating value estimates based on the difference (error) between predicted rewards and actually observed rewards over time.
- This difference is called the temporal difference error.

* TD Learning Process

- When an agent is in a particular state, it makes a prediction about the value of that state.
- The agent takes an action, receives a reward, and moves to the next state.
- It then makes a prediction about the value of the next state.
- The TD error is computed as the difference between the new prediction and the new prediction (adjusted by the received reward) and the old prediction.
- This error is used to update the value of the initial state, bringing it closer to the true value.

$$\text{TD Error} = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

Update Rule :

$$V(s_t) \leftarrow V(s_t) + \alpha \times \text{TD Error}$$

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

* Goal of TD Learning

- Estimate the value of states to improve decision-making.
- Imagine playing a game where each state corresponds to a situation with potential rewards.

(3)

→ You can take actions, which lead to new states and give one some points (rewards).

→ TD Learning helps understand how good each state is so that one can make better decisions and maximize one's total points (rewards) in the game.

* Why TD Learning

- (i) Fast Updates - Don't need to wait until an episode ends to update values. Updates occur after each action.
- (ii) Model Free - TD Learning doesn't require knowledge of the transition rules or environment model; it learns through experience alone.

* Examples of TD Learning

- ① Weather Prediction - You think today it will be sunny (current prediction). Suddenly, you see clouds forming (new info, a new state). You update your prediction to it might rain soon, based on the new state.

This update is like the TD error. You adjust your prediction because you learned something new.

- ② Delivery Robot - You have a delivery robot that needs to learn the best path to delivery packages from the lobby to various offices in a building.

The robot receives a reward for delivering the package successfully to the correct office. The robot's goal is to learn which paths lead to faster deliveries (higher rewards) and avoid paths that are slow or have obstacles (lower rewards).

* Steps in TD Learning

Step 1 : Initial Estimate

→ At the beginning, you have some guess about how good each state is. This is called the value of the state. Can start w/ random guesses initially

Step 2 : Experience a transition

→ Let the current state be s .

→ Take an action and move to a new state s'

→ Receive an ^{reward}~~action~~ r for this action

Step 3 : Make a prediction

→ Before taking the action, you had a prediction about the value of the state s .

→ After moving to state s' , you have a new prediction based on s and the reward received.

Step 4 : Calculate TD Error.

→ The TD error is the difference between what you predicted for state s and what you learned after moving to state s' .

* TD Learning vs. Monte Carlo Methods

TD Learning

1. One-step TD allows for updates at each step in an episode. - calculate reward diff. when moving from the current state to the next

$$2. V(s_t) \leftarrow V(s_t) +$$

$$\alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

Bootstrapping target

TD Error

Monte Carlo Methods

1. MC updates every state after analyzing the whole episode.

(roughly equivalent to an update for each state occurring after n steps)

$$V(s_t) \leftarrow V(s_t) +$$

$$\alpha [G_t - V(s_t)]$$

return

MC Error

(*) n-step Bootstrapping generalizes TD-algorithms

one-step TD Learning

$n=1$

Monte Carlo

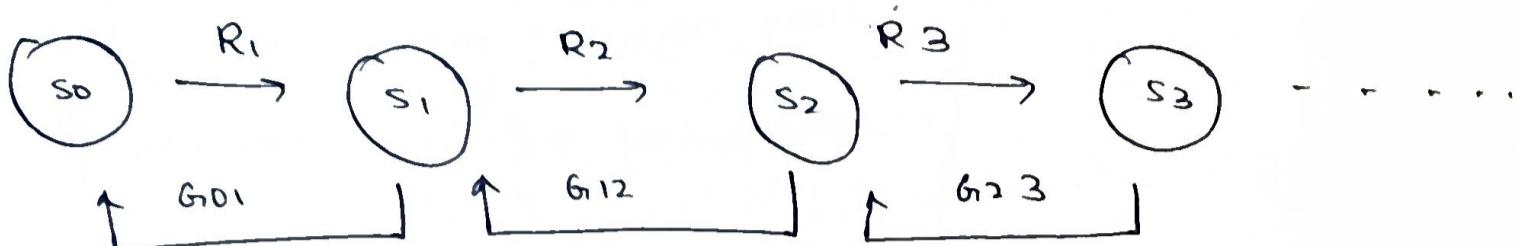
$n=\infty$

* n-step Bootstrapping

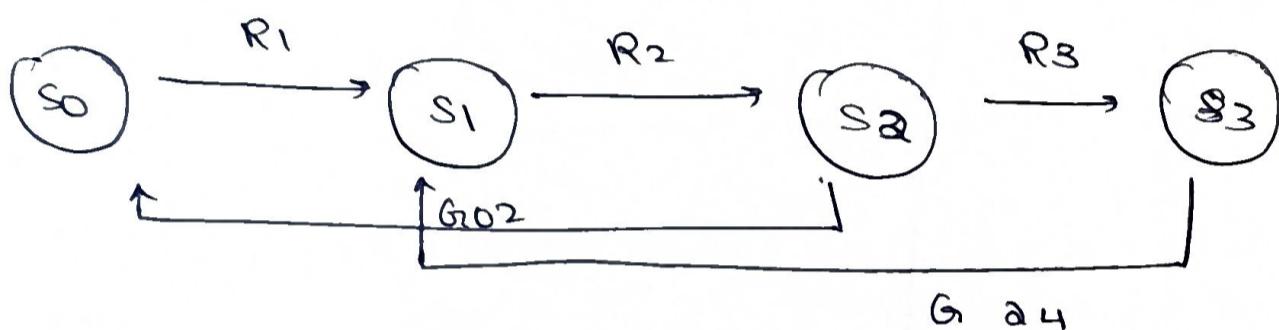
- The idea of 1-step TD learning can be extended to multiple steps.
- To do this, introduce n-step return, which calculates the accumulated discount reward, between the current state s_t and the future state at step s_{t+n} . In addition, it also adds state value at step s_{t+n} .

→ The update rule for n-step bootstrapping is

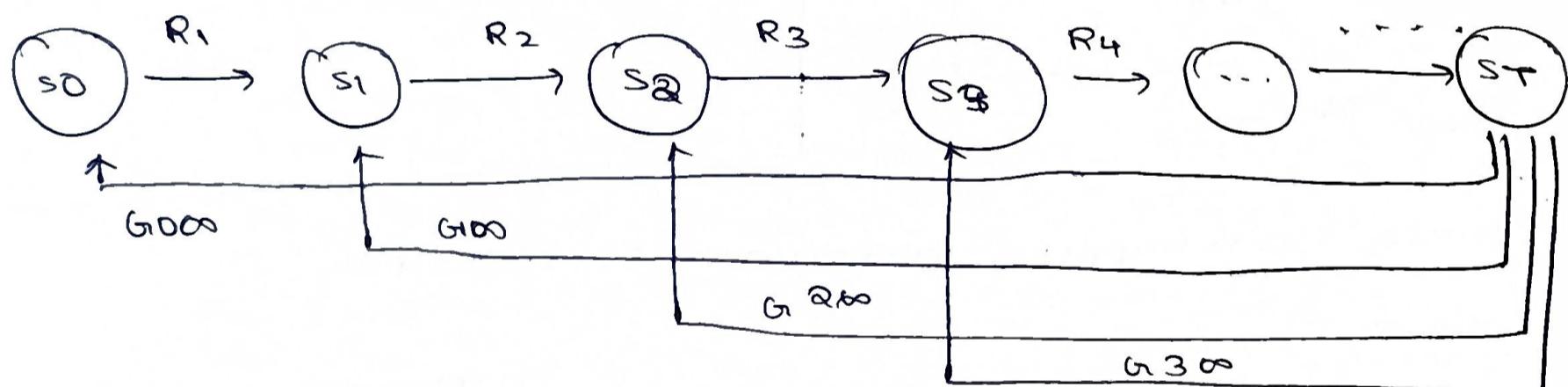
$$V(s_t) = V(s_t) + \alpha (G_{t+1:n} - V(s_t))$$



| 1-step TD



| 2-step TD



| ∞-step TD = Monte Carlo

example : in 3-step TD

- Beginning of episode is generated till state s_3
- s_0 is updated using 3-step return, which sums up rewards R_1, R_2, R_3 and the value of state s_3
- e.g. for $s_1 \rightarrow$ sum $R_2, R_3 \in R_4$ and s_4

7

Algorithm for n-step Bootstrapping

Input: policy π

Parameters: step size $\alpha \in (0,1]$, value in n

Initialize $v(s)$ arbitrarily $\forall s \in S$

Loop for each episode

Initialize & store s_0

$T \leftarrow \infty$

Loop for $t = 0, 1, 2, \dots$

If $t < T$ then

Take an action according to $\pi(\cdot | s_t)$

Observe and stored $R_{t+1} \& s_{t+1}$

If s_{t+1} is interminal $T \leftarrow t+1$

$\gamma \leftarrow t-n+1$

If $\gamma \geq 0$:

$$G_i \leftarrow \sum \gamma^{i-y-1} R_i$$

If $y+n < T$, $G_i \leftarrow G_i + \gamma^n v(s_{y+n})$

$$v(s_y) \leftarrow v(s_y) + \alpha [G_i - v(s_y)]$$

until $\gamma = T-1$

* On-Policy and Off-Policy Learning Methods

→ In RL, on-policy and off-policy methods refer to how an agent learns from the data it collects during interaction with an environment.

A. On-Policy Methods

Definition : On-policy methods learn the value of the same policy that the agent follows to make decisions. This means that the agent both explores and learns using the same policy.

Characteristics

- (i) Policy Improvement : The agent continuously improves the policy it is currently using to interact with the environment.
- (ii) Exploration and Exploitation Balance : Both managed by a single policy, that tries to balance trying new actions (exploration) with choosing the best-known action (exploitation).

e.g. ~~ε~~ ϵ -greedy policy

Example Algorithm - SARSA

Advantages - (i) more stable updates as they are based on the current behavior policy

Disadvantages - (i) slower convergence compared to off-policy methods
(ii) may not fully explore the environment if the exploration strategy is not effective.

B. OFF-Policy Methods

Definition - allows the agent to learn the value of an optimal policy independently of the actions it takes.

→ The agent can use one policy for exploration (called the behavior policy), while learning or improving another policy (often, an optimal policy, known as the target policy)

Characteristics -

Separate Learning and Exploration Policies - The agent gathers experience with one policy and learns with another, more optimal policy.

Example Algorithm - Q-Learning - agent uses the experience of its exploratory behaviour to update the value of the optimal policy.

Advantages (i) more flexible as it allows for off-policy exploration
 (ii) can learn from stored data (replay buffer), which is useful in environments with large state-action spaces

Disadvantages : (i) may be less stable due to the difference between the behavior and target policies.

(ii) May require more careful tuning of exploration strategies.

* Q-Learning

- Q-Learning is an RL technique that uses Q-values (or action-values) to guide an agent in selecting actions that yield the highest rewards over time.
- The primary goal is to build a Q-table that guides the agent towards optimal actions at each state.

Terms

- (i) Q-values : These values represent the quality of an action $Q(s,a)$ taken at a state s , indicating how rewarding the action is expected to be.
- $Q(s,a)$ is iteratively computed using the TD-update rule

- (ii) Reward : After each transition, the agent observes a reward for every action from the environment, and then transits to another state.

* Process of Q-Learning

- Initially, the agent explores the environment and updates the Q-Table. When the Q-Table is ready, the agent will start to exploit the environment and take better actions.
- It is an off-policy algorithm, meaning that it learns the optimal policy independently of the agent's current behavior policy, allowing it to refine its strategy.

→ Q-Learning estimates the value of Q at each time step, using the TD-update rule.

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

↑ reward
↑ max
TD - error

* Q-Learning Algorithm

1. Set the γ parameter
2. Set environment rewards in matrix
3. Initialize matrix Q as zero matrix

(rows denote the current state of the agent, columns denote the possible actions leading to the next state)

4. Select random initial ^{source} ~ _{state}

- Set initial state $s = \text{current state}$
- Select one action 'a' among all possible actions using the exploratory policy
- Take this action 'a', go to next state s'
- Observe reward r
- Get max Q-value to go to the next state based on all possible actions

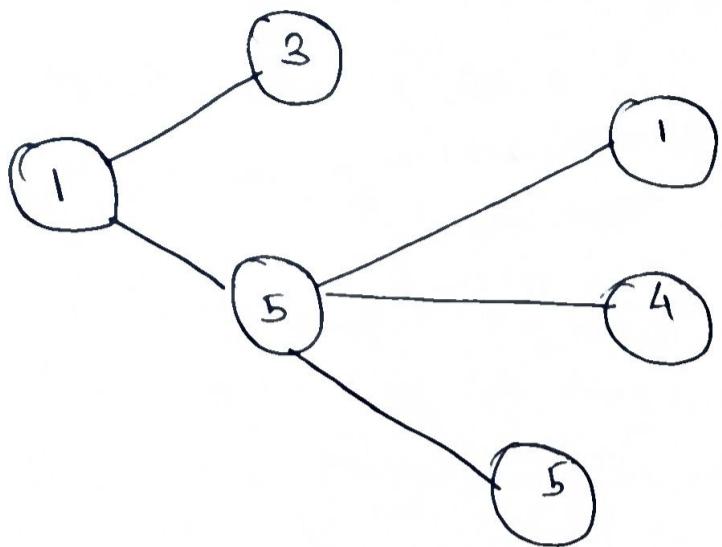
- b. Compute

$$Q(s,a) = Q(s,a) + \alpha (R + \gamma \max (Q(s',a') - Q(s,a)))$$

6. Repeat the above steps until current state = goal state.

Example

Consider the following graph:



$$\text{Let } \gamma = 0.8$$

The Q matrix is initialized to 0.

$$Q = \begin{bmatrix} 0 & 1 & 2 & 3 & 4 & 5 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

The R-matrix is

$$R = \begin{bmatrix} 0 & 1 & 2 & 3 & 4 & 5 \\ -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix}$$

From state 1, agent can

go to 3 or 5

go to 5 (say)

compute Q value

$$\Rightarrow Q(s,a) = R(s,a) + \gamma * \max [Q(\text{next state, all actions})]$$

$$\Rightarrow Q(1,5) = 100 + 0.8 * \max [Q(5,1), Q(5,4), Q(5,5)]$$

$$\Rightarrow Q(1,5) = 100 + 0.8 * \max (0, 0, 0)$$

$$Q(1,5) = 100 \rightarrow$$

Fill this value in

Q table

(B)

SARSA

- an RL algorithm that learns by interacting with its environment.
- an on-policy TD-learning algorithm. It updates the action-value function $Q(s,a)$, when starting from state s , taking action a , and following the current policy thereafter.
- The agent learns from its own experiences & refines its policy based on the observed rewards & state transitions.
- The update rule for SARSA is

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

Algorithm

1. Initialize Q-table - Fill with random values or zeroes, representing the estimated value of taking each action in each state.
2. Choose action - The agent chooses an action based on its current state using a policy, such as ϵ -greedy, which balances exploration & exploitation.
3. Take action - The agent takes the chosen action and observes the new state and the reward it receives from the environment.
4. Update Q-Table - The Q-table is updated using the SARSA update rule.
5. Repeat - The agent continues to iterate through the process, observing outcomes & updating its Q-table, till it converges.

* Advantages and Disadvantages of SARSA

Advantages

- more stable, less prone to oscillations compared to off-policy methods

Disadvantages

- can be sensitive to learning rate and exploration parameters
- can be slow to converge in environments w/ high variance or complex reward structures

* Planning and Learning with Tabular Methods

* Planning and Learning

Planning - any computational process that uses a model to create

- o. improve a policy

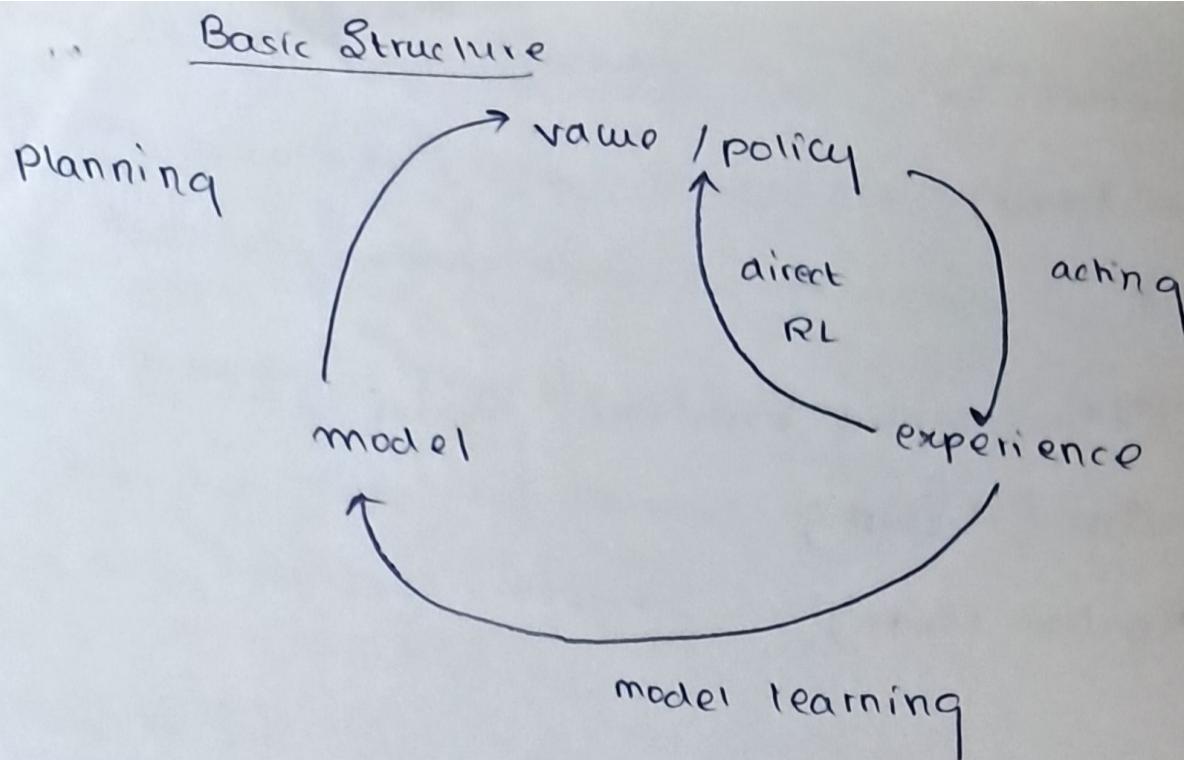
Learning - the acquisition of knowledge or skills through experience, study or being taught.

* Planning and Learning in RL

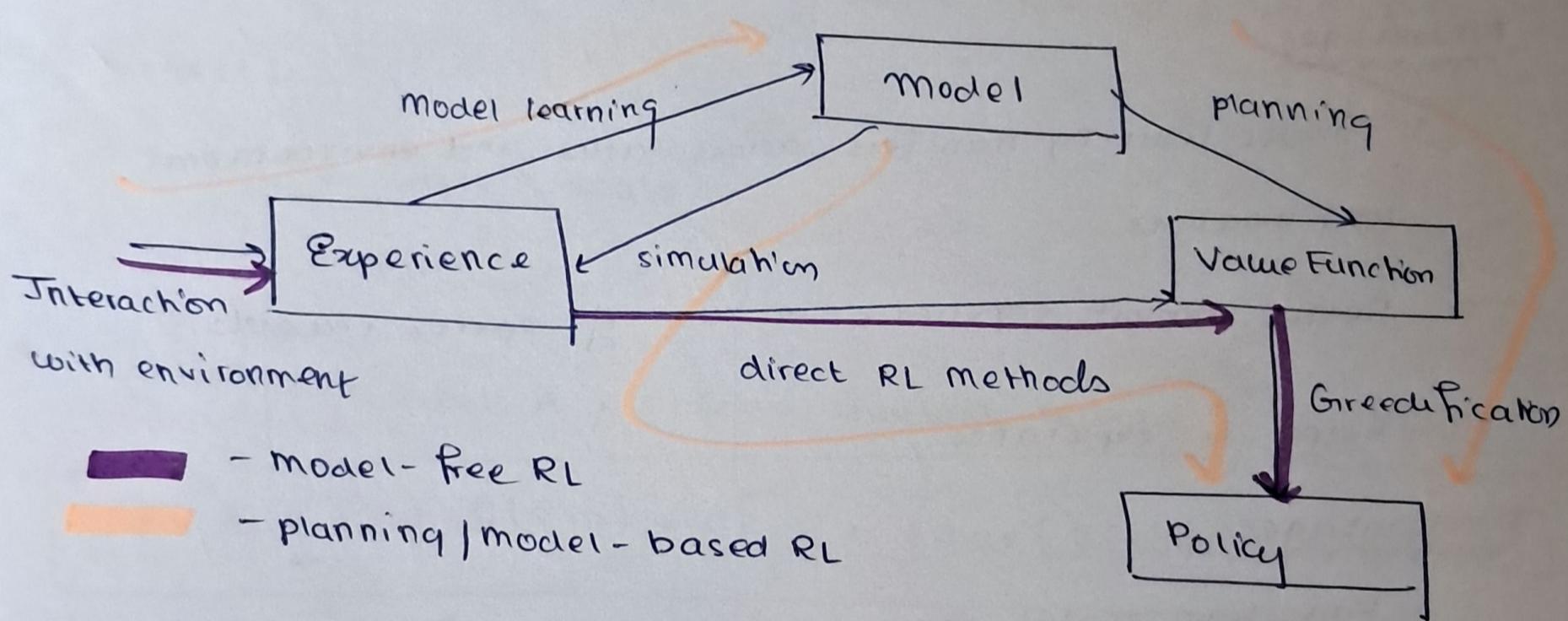
→ Planning involves the agent using a model of the environment to simulate interactions and update its policy w/o directly interacting with the environment

eg. Value Iteration, Policy Iteration

→ Learning involves the agent interacting with the environment to improve its knowledge about the reward structure / transition dynamics . eg. Q-Learning



Model-free vs. Model-based RL



* Features of Planning / Model-based RL

Combine real and simulated experience

1. If the model is unknown, learn the model
2. Learn value functions using both real experience & the model
3. Learning value functions online using model-based lookahead search

i.e let the real experience be:

$$s' \sim T(s'|s,a)$$

$$R = r(s,a)$$

and the simulated experience, sampled from the model is

$$s' \sim T_\eta(s'|s,A)$$

$$R = R_\eta(R|s,A)$$

* Advantages and Disadvantages of model-based RL

Advantages

1. Model learning transfers across tasks and environment configurations
2. Better exploits experience in case of sparse rewards
3. Helps in exploration

Disadvantage

1. First learn model, then construct a value function: Two sources of approximation error

* Dyna-Q

→ In dyna-Q, learning and planning are accomplished by exactly the same algorithm, operating on real experience for learning and on simulated experience for planning.

Key Concepts

- ① Direct RL: Involves the agent updating its Q-values directly from real interactions with the environment.

(2) Model learning - The agent learns a model of the environment by storing the transitions it observes. This model can be deterministic or stochastic. (17)

(3) Planning - Involves updating Q-values based on simulated experience generated from the learned model. This allows the agent to improve its policy without requiring additional interactions with the environment.

Algorithm

Initialize $Q(s,a)$ and $\text{Model}(s,a)$ for $\forall s \in S$ and $a \in A(s)$

Do forever:

(a) $s \leftarrow$ current state

(b) $a \leftarrow \epsilon_t\text{-greedy}(s, Q)$

(c) Execute action A , observe resultant reward R and state s'

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

↳ Direct RL

(e) $\text{Model}(s,a) \leftarrow R, s'$

↳ Model learning

↳ Planning

(f) Repeat n times:

$s \leftarrow$ random previously observed state

$a \leftarrow$ random action previously taken in s

$R, s' \leftarrow \text{Model}(s, a)$

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

* Prioritized Sweeping?

- Simulating with random transitions might be very heavy in some applications.
- Instead of performing uniform or random updates (as in Dyna-Q), Prioritized Sweeping selectively updates state-action pairs that are more likely to lead to significant changes in the value function
- A queue is maintained of every state-action pair whose estimated value would change non-trivially if updated, prioritized by the size of the change.
- Prioritized sweeping has been found to dramatically increase the speed at which optimal solutions are found in maze tasks, often by a factor of 5 to 10.
- The priority is calculated based on the magnitude of the temporal difference (TD) - error

$$\text{priority } (s', a') = |R + \gamma \max_a Q(s', a) - Q(s, a)|$$

- If the TD error is large, it indicates that the value estimate for (s', a') is significantly outdated, so it should be prioritized for updating

Steps in Prioritized Sweeping

1. Compute Priority - After observing a transition, and updating the value function for a state-action pair (s, a) , the algorithm computes the priority for predecessor states that can transition into s via some action a'
2. Push into priority queue - The calculated priority is then pushed into the priority queue. The larger the TD-error, the higher the priority, meaning this state-action pair will be updated sooner
3. Prioritized Updates - The planning loop then performs updates on state-action pairs with the highest priority first

Algorithm

Initialize $Q(s, a)$, Model (s, a) for $\forall s, a$ and PQqueue to empty
Do forever:

- (a) $s \leftarrow$ current state
- (b) $a \leftarrow$ policy (s, Q)
- (c) Execute action a , observe reward R and state s'
- (d) Model $(s, a) \leftarrow R, s'$
- (e) $P \leftarrow |R + \gamma \max_a Q(s', a) - Q(s, a)|$
- (f) IF $P > 0$, then insert s, a into PQqueue with priority P
- (g) Repeat n times, while PQqueue is not empty :

$s, a \leftarrow \text{first}(\text{PQueue})$

$r, s' \leftarrow \text{model}(s, a)$

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

Repeat $\forall \bar{s}, \bar{a}$ predicted to lead to s :

$\bar{r} \leftarrow \text{predicted reward for } \bar{s}, \bar{a}, s$

$$P \leftarrow |\bar{r} + \gamma \max_a Q(s, a) - Q(\bar{s}, \bar{a})|$$

$\text{if } P > 0 \text{ then insert } \bar{s}, \bar{a} \text{ into PQueue with priority } P$

Disadvantage of Prioritized Sweeping

→ It uses expected updates, which in stochastic environments may waste lots of computation on low-probability transitions.