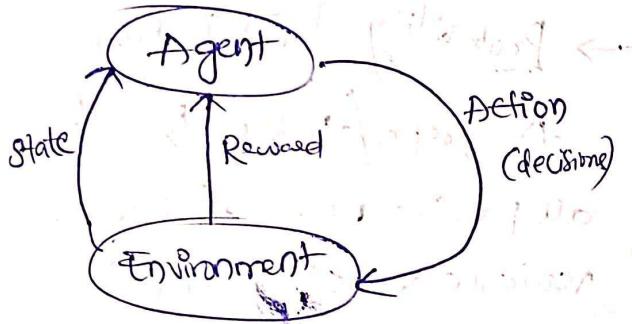


Reinforcement Learning -

It's a type of ML, where an agent takes action based on the changes occurring in the environment (feedback). It stores the reward received by performing the action in state stage. Using stored data it takes next step & continues for further actions.

Here the model learns by itself.



- The goal is to achieve the best result.
- RL is totally different from supervised & unsupervised learning.
- RL is an area of ML. It is about taking suitable action to maximize the reward in a particular situation.
- In RL, RL agent decides what to perform according to the situations for a given task. RL agent learns by itself from its experience.

Key Components

1) Agent: - The learner & decision maker tries to learn the best actions to take in an environment.

2) Environment:

This is the place where an agent will interact with. It can be anything like game, real world tasks etc.

3) Action:

These are the moves taken by an agent in the environment.

4) State:

State is a situation returned by the environment after each action is taken by the agent.

5) Reward:

The feedback returned to agent from the environment for its action.

6) Policy: - The strategy or rule that the agent follows to decide the actions to take in different states.

→ There are mainly two types of RL.

1) Model RL

2) Model free RL

→ In model based RL The agent tries to build the model and has access to a model of the environment.

→ Once the agent has the model, agent can plan ahead and choose actions.

→ In model free RL the agent doesn't build the model of environment. The agent learns without a model of the environment.

→ Basics of probability of linear algebra: To understand RL one should have basic of Probability Basics in RL:

Probability: It is the measure of the likelihood of an event occurring.

It lies b/w 0 & 1.

In RL Probabilities are used to represent uncertainty such as likelihood of transitioning from one state to another state etc.

$$P(E) = \frac{\text{no. of favourable outcomes}}{\text{Total no. of outcomes}}$$

↓
Event

Ex: Coin a toss.

$$P(H) = \frac{1}{2} \begin{matrix} (\text{Favourable}) \\ \text{Head} \end{matrix}$$

Random Variable:

Variable that take different values randomly.

In RL, RV represent states, actions, reward and other quantities that are subjected to uncertainty (no perfect information)

→ Probability Distribution:

A graph/table that shows all possible values of a random variable & their probabilities.

→ Expected Value:

The avg value of a RV.

Probability distribution Ex:

Coin a Toss.

$$P(H) = \frac{1}{2} \quad (\because \text{All values are shown})$$

$$P(T) = \frac{1}{2}$$

→ Conditional Probability:

The prob. of an event occurring given that another event has occurred.

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

Linear Algebra in RL!

LA is the branch of mathematics concerning linear equation such as

$$a_1x_1 + a_2x_2 + \dots + a_nx_n = b$$

Eigen values & vectors:

These are used in algorithms like PCA (principal component analysis) to extract important features from states.

Linear algebra helps in managing & processing data when the environment is complex & when we have lot of information i.e. (large state space & actions).

In RL LA is used to represent the relationship b/w variables to make it easier to analyse & minimize the learning process.

Linear algebra includes:

Vectors: These are mathematical obj's that represent both magnitude & direction.

In RL things like states, actions & rewards can be represented as vectors.

Matrices:

These are used to represent transition probabilities, rewards & value function.

Matrix operations:

Such as Add, sub, mul, & inversion. Are used in RL to update value functions, calculate state-action values

→ Both probability & linear algebra are important for RL bcoz probability helps us model the uncertainty in RL environment whereas LA helps us making decisions by representing value fun. of policy.

→ Even both help us optimize the performance of RL agent.

→ Definition of Stochastic Multi Armed Bandit:

A stochastic multi armed bandit (MAB) is a classic problem in RL that involves making decisions in an uncertain environment.

MAB is a simplified decision making problem where we have multiple slot machines each with different unknown reward probabilities.

Arm refers to an option or choice that an agent can select.

→ In simple SMAB is choosing the best machine with the given multiple slot machines (arms).

→ The aim of SMAB is to find the optimal arm to maximize the cumulative reward (total reward) or sum of all the rewards.

→ Imagine you are in a casino with several slot machines (arms). Each machine has a different prob. of paying out money when you will pull the lever.

Initially we don't know the probabilities of machines in advance.

→ Our goal is to figure out which machine is the best and gives higher reward by pulling the lever.

→ Working

Initially there are multiple machines each provides different rewards.

We need to find the best one by exploration and exploitation.

→ Exploration defines trying different machines to learn which is best.

→ Exploitation focus on the machine that already known to give the best rewards.

⇒ Definition of Regret :

In RL regret is a measure of how much worse ^(worst) your agent's performance is compared to the best possible performance it could have achieved.

Regret is the difference b/w the optimal reward an agent could have received and the actual reward it received.

$$\text{Regret} = \frac{\text{Optimal Reward} - \text{Actual Reward}}{\text{Optimal Reward}}$$

→ Regret is a crucial concept in RL bcoz it helps agents evaluate their decisions and learns from their mistakes.

By minimizing regret, agent can improve their performance & make better decisions over time.

Ex:

Suppose an agent is playing a game where it receives

seconds for completing levels. The reward for completing 1 level is 100 points. However, the agent only receives 80 points.

The regret would be

$$\text{Regret} = \text{Optimal Reward} - \text{Actual Reward}$$
$$100 - 80$$

$$\boxed{R = 20.}$$

Here agent missed out 20 points.

→ 2 types of regret

- 1) Internal
- 2) External

It measures the agent's actual reward if the reward it would have received, is called external regret. If it followed best fixed action.

Internal regret:

It measures the difference b/w the agent's actual reward & the reward it would have received if it had followed the best sequence of action.

→ Achieving Sublinear regret (ASR):

ASR in RL refers that the agent interacts with the environment over time, the amount of regret it experiences grows more slowly than the total no. of actions it takes.

Sub linear regret means the agent's regret mistakes decreases over time.

→ ASR is important bcz it ensures that the agent learns efficiently & make good decisions most of the time.

→ ASR minimizes the difference b/w total reward obtained by an agent & the total reward of an optimal strategy.

→ Strategies / Types of ASR:

- 1) Explore-Exploit trade off / Strategy

The agent needs to balance trying out new actions by exploring and sticking with the actions that have already worked well by exploiting.

- 2) Upper Confidence Bound (UCB):

Assigning confidence bounds to each actions estimated value.

Ex: Consider 2 actions, A & B. If agent tried many times action A, & B tried only few times then UCB choose action B to reduce uncertainty.

3) Thompson Sampling \Rightarrow UCB Algorithm:

This strategy uses probability distribution to choose actions.

It aims to maximize rewards by balancing exploration & exploitation.

If action A is mostly likely to give the highest reward, the agent will pick action A, but may still explore actions based on their probabilities.

4) optimism in the face of uncertainty:

using optimistic initial values to encourage exploration, i.e. be optimistic abt unknown actions & try them.

In RL, the idea of optimism is often used when agent isn't sure which actions will lead to the best reward.

\rightarrow ASR achieves good decision making & improved performance.

upper confidence bound.

- It's an algo. used in RL to Balance Exploration & Exploitation.

- It is a way to make decisions in situations where you don't know the best option.

Ex: choosing best machine among the multiple machines.

UCB working

UCB focuses on maximizing reward over-time by considering both how good an option seems and how much you need to explore that option.

1) Initialize the estimated reward and confidence bounds for each action.

2) Choose the actions with the highest UCB.

3) Update the estimated reward & CB based on the observed reward.

& repeat step 2, 3

\rightarrow UCB gives higher priority to options that have been tried less

Types of UCB

1) UCB1 (classic UCB) :-

This is the basic version of UCB. It is the original UCB algo. that balances explore & exploitation using UCB formula.

formula

$$UCB(A) = \text{avg reward of actions } A + C \times \sqrt{\frac{\ln(t)}{N(a)}}$$

where

$UCB(a)$ = Upper confidence bound for action 'a'.

avg reward → taking actions.

C = parameter controlling b/w explor & exploit.

$\ln(t)$ = natural algorithm.

$N(a)$ = no. of times action has been chosen

2) UCB Tuned :-

This is a refined version of UCB1.

It adjusts the constant factor in the confidence bound to handle reward distribution for action.

It outperforms UCB1 in many environments.

3) KL-UCB

Kullback-Leibler UCB,

It focuses on minimizing the KL divergence b/w the estimated & true reward distribution.

4) MOS : (Minimax optimistic Sampling Strategy) :-

It is the extension of UCB.

It provides more conservative approach by minimizing.

→ KL-UCB :-

- It is the another version of the UCB algo, but it uses more refined way to balance exploration & exploitation.

→ The key difference with regular UCB is that KL-UCB uses a more probabilistic approach to measure uncertainty.

→ Instead of just considering how many times an action has been selected, it looks at how far the current distribution of rewards for each action is from what you expected and adjust the exploration accordingly.

→ KL-UCB uses KL-divergence to calculate the UCB for each action.

- It helps the agent make decisions by considering both estimated value of an action & uncertainty of information.
- KL-UCB tries to make smottee decision using a more mathematical & detailed understanding of the uncertainty in the reward.

Working:

- 1) The agent start with an initial estimate of the reward distribution for each action.
- 2) The agent calculates the KL divergence b/w estimated reward distribution and a uniform distribution.
- 3) The agent uses KL divergence to calculate a confidence bound for each action.
- 4) The agent chooses the action with highest UCB.
- 5) The agent updates its estimate of the reward for the chosen action based on the observed reward.

6) The agent repeats steps 2-5.

→ KL-divergence measures the difference b/w 2 probability distributions.

→ confidence bound estimates the uncertainty of the reward distribution

→ UCB is the maximum possible reward for an action

KL-UCB Formula

$$\text{UCB}(a) = Q(a) + \sqrt{\frac{2 \log(t)}{N(a)} \text{KL}(P_{\text{eq}})}$$

where

$\text{UCB}(a)$ = UCB of an action a .

$Q(a)$ = estimated reward for action a :

P is estimated reward distribution

q is uniform distribution.

$\text{KL}(P_{\text{eq}})$ is KL divergence b/w

P & q

it is current time step.

$N(a)$ is no. of times our action has been selected / chosen.

Thompson Sampling :-

It is a method used in RL to help an agent make decisions based on uncertainty.

- It is an algo. that helps agent to choose actions in an uncertain environment.
- It involves in maintaining a probability distribution over possible actions to decide next action.

→ There are 2 methods.

1) Bernoulli Thompson Sampling

2) Gaussian " "

→ BTS is used for binary outcomes (Success/Failure)

Eg: Coin a Toss, if Head Success else fail

GTS is used for continuous outcomes.

Working:-

- Initialize a probability distribution over the possible rewards for each action.
- Sample a reward from the probability distribution for each action.
- Choose the action with highest sampled reward.
- Update the PD based on observed reward.

Formula:-

$$\alpha(a) \leftarrow \alpha(a) + \gamma$$

$$\beta(a) \leftarrow \beta(a) + 1 - \gamma$$

$\alpha(a)$ & $\beta(a)$ are the shape parameters of the β -distribution for action (a) .

$\gamma \Rightarrow$ observed reward.