

## Unit-V Reinforcement Learning

Reinforcement learning (RL) is a machine learning technique that focuses on how AI agents should take actions in an environment to maximize the total reward. The training is done in real time with continuous feedback to maximize the possibility of being rewarded.

Reinforcement learning lets a machine learn from its mistakes, similar to how humans do. It's a type of machine learning in which the machine learns to solve a problem using trial and error. Also, the machine learns from its actions, unlike supervised learning, where historical data plays a critical role.

The AI system that undergoes the learning process is called the agent or the learner. The learning system explores and observes the environment around it, just like us. If the agent performs the right action, it receives positive feedback or a positive reward. If it takes an adverse action, it receives negative feedback or a negative reward.

### **Notable characteristics of reinforcement learning (RL) are:**

Time plays a critical role in RL problems.

The agent's decision-making is sequential.

There isn't a supervisor, and the agent isn't given any instructions. There are only rewards.

The agent's actions directly affect the subsequent data it receives.

The agent is rewarded (positive or negative) for each action.

The best solution to a problem is decided based on the maximum reward.

### **Terminologies used in reinforcement learning**

**Agent:** The AI system that undergoes the learning process. Also called the learner or decision-maker. The algorithm is the agent.

**Action:** The set of all possible moves an agent can make.

**Environment:** The world through which the agent moves and receives feedback. The environment takes the agent's current state and action as input and then outputs the reward and the next state.

**State:** An immediate situation in which the agent finds itself. It can be a specific moment or position in the environment. It can also be a current as well as a future situation. In simple words, it's the agent's state in the environment.

**Reward:** For every action made, the agent receives a reward from the environment. A reward could be positive or negative, depending on the action.

**Policy:** The strategy the agent uses to determine the next action based on the current state. In other words, it maps states to actions so that the agent can choose the action with the highest reward.

**Model:** The agent's view of the environment. It maps the state-action pairs to the probability distributions over states. However, not every RL agent uses a model of its environment.

**Value function:** In simple terms, the value function represents how favorable a state is for the agent. The state's value represents the long-term reward the agent will receive starting from that particular state to executing a specific policy.

**Discount factor:** Discount factor ( $\gamma$ ) determines how much the agent cares about rewards in the distant future when compared to those in the immediate future. It's a value between zero and one. If the discount factor equals 0, the agent will only learn about actions that produce immediate rewards. If it's equal to 1, the agent will evaluate its actions based on the sum of its future rewards.

**Dynamic programming (DP):** An algorithmic technique used to solve an optimization problem by breaking it down into subproblems. It follows the concept that the optimal solution to the overall problem depends on the optimal solution to its subproblems.

## Types of reinforcement learning

There are two types of reinforcement learning methods: positive reinforcement and negative reinforcement.

### Positive reinforcement

Positive reinforcement learning is the process of encouraging or adding something when an expected behavior pattern is exhibited to increase the likelihood of the same behavior being repeated.

For example, if a child passes a test with impressive grades, they can be positively reinforced with an ice cream cone.

### Negative reinforcement

Negative reinforcement involves increasing the chances of specific behavior to occur again by removing the negative condition.

For example, if a child fails a test, they can be negatively reinforced by taking away their video games. This is not precisely punishing the child for failing, but removing a negative condition (in this case, video games) that might have caused the kid to fail the test.

## Elements of reinforcement learning

Apart from the agent and the environment, there are four critical elements in reinforcement learning: policy, reward signal, value function, and model.

### 1. Policy

The policy is the strategy the agent uses to determine the following action based on the current state. It's one of the critical elements of reinforcement learning and can single-handedly define the agent's behavior.

A policy maps the perceived states of the environment to the actions taken on those particular states. It can be deterministic or stochastic and can also be a simple function or a lookup table.

## **2. Reward signal**

At each state, the agent receives an immediate signal from the environment called the reward signal or simply reward. As mentioned earlier, rewards can be positive or negative, depending on the agent's actions. The reward signal can also force the agent to change the policy. For example, if the agent's actions lead to negative rewards, the agent will be forced to change the policy for the sake of its total reward.

## **3. Value function**

Value function gives information about how favorable specific actions are and how much reward the agent can expect. Simply put, the value function determines how good a state is for the agent to be in. The value function depends on the agent's policy and the reward, and its goal is to estimate values to achieve more rewards.

## **4. Model**

The model mimics the behavior of the environment. Using a model, you can make inferences about the environment and how it'll behave. For example, if a state and an action are provided, you can use a model to predict the next state and reward.

Since the model lets you consider all the future situations before experiencing them, you can use it for planning. The approach used for solving reinforcement learning problems with the model's help is called model-based reinforcement learning. On the other hand, if you try solving RL problems without using a model, it's called model-free reinforcement learning.

While model-based learning tries to choose the optimal policy based on the learned model, model-free learning demands the agent learn from trial-and-error experience. Statistically, model-free methods are less efficient than model-based methods.

## **Discounting in RL**

Discounting in reinforcement learning is a technique used to give more importance to immediate rewards compared to future rewards. It involves applying a discount factor to future rewards, which reduces their relative value as they are temporally further away.

The purpose of discounting is to account for the inherent uncertainty and delayed consequences in reinforcement learning problems. By assigning a lower weight to future rewards, the agent becomes more inclined to prioritize immediate rewards, leading to faster convergence and better decision-making in dynamic environments.

Mathematically, the discount factor (usually denoted as  $\gamma$ , gamma) is a value between 0 and 1. It determines the extent of discounting for future rewards. A discount factor of 1 indicates no discounting, meaning future rewards are valued equally to immediate rewards. A discount factor of 0 means that only immediate rewards are considered, and future rewards have no influence on decision-making.

The discounted cumulative reward, also known as the discounted return, is calculated as the sum of discounted rewards over time. It is defined by the formula:

$$R(t) = r(t) + \gamma * r(t+1) + \gamma^2 * r(t+2) + \gamma^3 * r(t+3) + \dots$$

where  $r(t)$  represents the reward at time  $t$ , and  $\gamma$  is the discount factor.

By applying discounting, the agent learns to optimize its actions in a way that balances immediate rewards with the long-term cumulative reward. The choice of the discount factor depends on the specific problem and the trade-off between immediate rewards and future rewards in the given environment. A higher discount factor values future rewards more, leading to a more far-sighted decision-making process, while a lower discount factor prioritizes immediate rewards.

Discounting in reinforcement learning helps in addressing the challenge of delayed rewards and assists the agent in making effective decisions by considering both short-term and long-term consequences.

## Working of reinforcement learning

Simply put, reinforcement learning is an agent's quest to maximize the reward it receives. There's no human to supervise the learning process, and the agent makes sequential decisions.

Unlike supervised learning, reinforcement learning doesn't demand you to label data or correct suboptimal actions. Instead, the goal is to find a balance between exploration and exploitation.

Exploration is when the agent learns by leaving its comfort zone, and doing so might put its reward at stake. Exploration is often challenging and is like entering uncharted territory. Think of it as trying a restaurant you've never been to. In the best-case scenario, you might end up discovering a new favorite restaurant and giving your taste buds a treat. In the worst-case scenario, you might end up sick due to improperly cooked food.

Exploitation is when the agent stays in its comfort zone and exploits the currently available knowledge. It's risk-free as there's no chance of attracting a penalty and the agent keeps repeating the same thing. It's like visiting your favourite restaurant every day and not being open to new experiences. Of course, it's a safe choice, but there might be a better restaurant out there.

Reinforcement learning is a trade-off between exploration and exploitation. RL algorithms can be made to both explore and exploit at varying degrees.

Reinforcement learning is an iterative process. The agent starts with no hint about the rewards it can expect from specific state-action pairs. It learns as it goes through these states multiple times and eventually becomes adept. In short, the agent starts as a noob and slowly becomes a pro.

## Markov Decision Process

A Markov Decision Process (MDP) is a mathematical framework used to model decision-making problems in a stochastic environment. It provides a formal representation of a sequential decision-making process where an agent interacts with an environment.

In an MDP, the environment is modeled as a set of states, actions, rewards, and transition probabilities. The key components of an MDP are as follows:

States (S): The set of possible states that the environment can be in. The agent makes decisions based on the current state.

Actions (A): The set of actions that the agent can choose from at each state. The agent selects an action to transition from the current state to the next state.

Transition Probabilities (P): The probabilities that define the likelihood of transitioning from one state to another after taking a particular action. These probabilities are represented by the transition function  $P(s, a, s')$ , where  $s$  is the current state,  $a$  is the action, and  $s'$  is the next state.

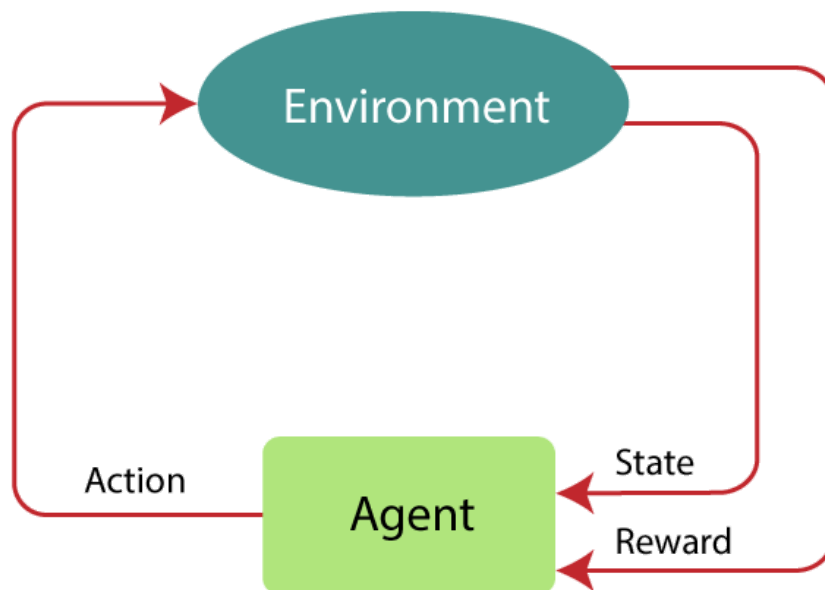
Rewards (R): The immediate rewards or penalties associated with taking specific actions in certain states. Rewards are usually represented by the reward function  $R(s, a, s')$  that provides the reward obtained when transitioning from state  $s$  to state  $s'$  after taking action  $a$ .

Discount Factor ( $\gamma$ ): A value between 0 and 1 that determines the importance of future rewards. It allows the agent to balance immediate rewards against long-term cumulative rewards. A higher discount factor places more emphasis on future rewards.

The goal in an MDP is to find an optimal policy that maximizes the expected cumulative reward over time. A policy is a mapping from states to actions, specifying the action the agent should take in each state. The optimal policy is the one that maximizes the expected cumulative reward.

Solving an MDP involves finding the optimal value function or Q-function, which represents the expected cumulative reward for taking a specific action in a specific state. Value iteration or methods like Q-learning can be used to estimate the optimal value function and derive the optimal policy.

MDPs find applications in various fields, including robotics, automated control systems, economics, healthcare, and more. They provide a mathematical framework for modeling and solving decision-making problems in uncertain and dynamic environments.



## Qlearning

Q-learning is a popular algorithm in reinforcement learning (RL) used for learning optimal policies in Markov Decision Processes (MDPs) without requiring a model of the environment. It is a model-free, off-policy algorithm that uses a value-based approach to iteratively update action-value estimates.

Here are the key steps of the Q-learning algorithm:

### Initialization:

Initialize the Q-values, denoted as  $Q(s, a)$ , for all state-action pairs  $(s, a)$  in the MDP. Q-values can be initialized arbitrarily or to some default value.

### Action Selection:

Choose an action,  $a$ , based on the current state,  $s$ , using an exploration-exploitation strategy. Common strategies include epsilon-greedy or softmax exploration to balance exploration of new actions and exploitation of learned knowledge.

Perform Action and Observe Reward and Next State:

Take action,  $a$ , in the environment.

Observe the reward,  $r$ , received from the environment.

Transition to the next state,  $s'$ , based on the dynamics of the environment.

### Update Q-Value:

Update the Q-value for the current state-action pair based on the observed reward and the maximum Q-value of the next state.

Use the Q-value update equation:

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

where  $\alpha$  is the learning rate ( $0 \leq \alpha \leq 1$ ) and  $\gamma$  is the discount factor ( $0 \leq \gamma \leq 1$ ) that represents the importance of future rewards.

Repeat Steps 2-4:

Repeat Steps 2 to 4 until a termination condition is met, such as reaching a maximum number of iterations or achieving convergence of Q-values.

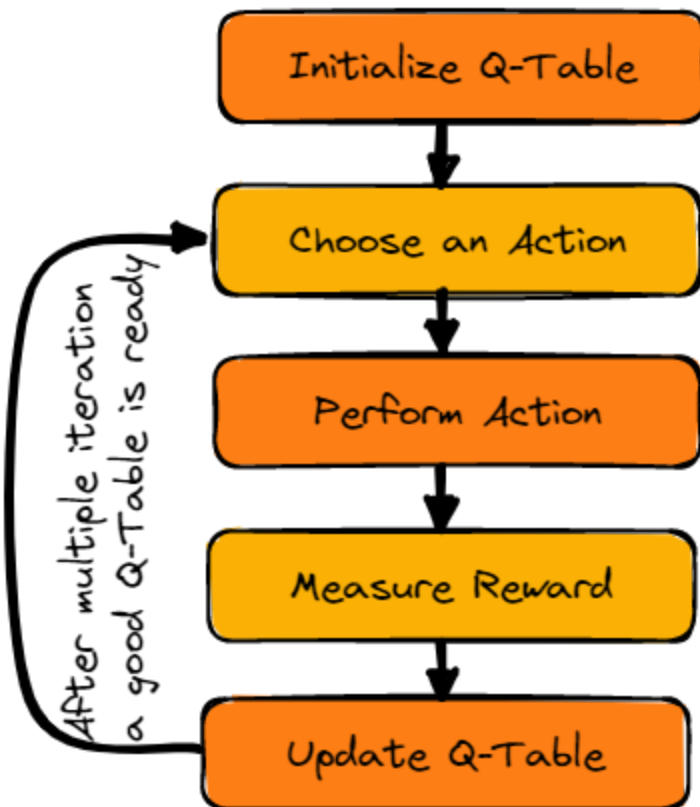
### Policy Extraction:

After learning the Q-values, the optimal policy can be derived by selecting the action with the highest Q-value for each state:  $\pi(s) = \operatorname{argmax}[Q(s, a)]$ .

### Repeat Episodes:

Repeat the entire process for multiple episodes, where each episode represents an interaction with the environment from an initial state until a terminal state is reached.

Q-learning updates the Q-values iteratively based on the observed rewards and the maximum Q-value of the next state. Through this process, the Q-values gradually converge to the optimal values, enabling the agent to make optimal decisions.



Note that Q-learning is an off-policy algorithm, meaning it learns the Q-values independently of the policy being followed. It updates the Q-values based on the maximum Q-value of the next state, regardless of the action actually chosen. This allows Q-learning to learn the optimal policy while still exploring the environment.

Q-learning has been widely used in various RL applications, including robotic control, game playing, autonomous navigation, and more.

# Applications of reinforcement learning

- Business strategy planning
- Aircraft control and robot motion control
- Industrial automation
- Data processing
- Augmented NLP
- Recommendation systems
- Bidding and advertising
- Traffic light control

## ML applications in Text classification

Machine learning (ML) has various applications in text classification, where the goal is to automatically categorize or assign predefined labels to text documents based on their content. Text classification plays a significant role in a wide range of areas, including:

- **Sentiment Analysis:** ML algorithms can classify text documents (such as customer reviews, social media posts, or news articles) into positive, negative, or neutral sentiment categories. This enables businesses to analyze public opinion and gauge customer sentiment towards products, services, or brands.
- **Spam Filtering:** ML models can be trained to distinguish between legitimate emails and spam emails by analyzing the content, structure, and patterns within the text. This helps in identifying and filtering out unwanted or malicious messages.
- **Document Categorization:** ML techniques can automatically classify text documents into specific categories or topics, such as news articles into sports, politics, entertainment, or technology. This aids in organizing and indexing large collections of documents for efficient retrieval.
- **Language Detection:** ML algorithms can determine the language of a given text document based on its linguistic characteristics. This is useful in various applications like multilingual customer support, content localization, and language-specific analysis.
- **Intent Recognition:** ML models can classify user queries or natural language inputs into specific intents or actions, such as determining if a user wants to ask a question, make a reservation, or seek assistance. This is commonly used in chatbots, virtual assistants, and customer support systems.
- **Topic Modeling:** ML algorithms, such as Latent Dirichlet Allocation (LDA), can automatically discover latent topics within a collection of documents. This helps in uncovering hidden themes, identifying common patterns, and organizing large textual datasets.
- **Textual Fraud Detection:** ML models can learn patterns of fraudulent or malicious activities by analyzing text data associated with financial transactions, insurance claims, or online user behavior. This helps in detecting and preventing fraudulent activities.

These are just a few examples of how ML is applied in text classification. ML techniques like Naive Bayes, Support Vector Machines (SVM), Random Forests, and Neural Networks are commonly used in these applications to learn patterns and make predictions based on textual data.



## ML applications in Image classification

Machine learning (ML) has numerous applications in image classification, where the objective is to automatically assign predefined labels or categories to images based on their visual content. Image classification finds applications in various fields, including:

- **Object Recognition:** ML algorithms can be trained to recognize and classify specific objects or entities within images, such as identifying different species of animals, types of vehicles, or recognizing common objects like chairs, tables, or buildings.
- **Facial Recognition:** ML models can be used to detect and recognize faces in images, allowing for applications like identity verification, access control, or sentiment analysis based on facial expressions.
- **Medical Image Analysis:** ML algorithms can aid in the analysis of medical images, such as X-rays, MRI scans, or histopathology slides. They can assist in automated diagnosis, identifying abnormalities, or detecting specific medical conditions based on image patterns.
- **Image-based Document Classification:** ML models can classify documents based on their visual content, such as distinguishing between handwritten and printed text, identifying different document types (invoices, forms, passports), or extracting specific information from documents.
- **Fine-grained Image Classification:** ML algorithms can classify images into fine-grained categories that require detailed discrimination, such as differentiating between different bird species, dog breeds, or flower varieties.
- **Visual Search:** ML-powered image classification enables visual search, where users can search for similar images based on a query image. This finds applications in e-commerce, art collections, or image databases where users can search for visually similar items.
- **Quality Control and Defect Detection:** ML models can be used to analyze images in manufacturing or production processes to detect defects, identify quality issues, or perform automated visual inspections.
- **Autonomous Vehicles:** ML algorithms play a crucial role in enabling object detection and classification for autonomous vehicles. They help in identifying and recognizing objects such as pedestrians, traffic signs, or other vehicles to make real-time decisions.

These are just a few examples of how ML is applied in image classification. ML techniques such as Convolutional Neural Networks (CNNs), Transfer Learning, and Deep Learning frameworks like TensorFlow or PyTorch are commonly used to build effective image classification models by learning patterns and features from images.

## ML applications in speech recognition

Machine learning (ML) has significant applications in speech recognition, which involves converting spoken language into written text or interpreting and understanding spoken commands. ML techniques enable the development of accurate and efficient speech recognition systems. Some of the key applications of ML in speech recognition include:

- **Speech-to-Text Transcription:** ML algorithms are used to convert spoken language into written text. This application finds use in transcription services, voice assistants, call center analytics, and accessibility tools for individuals with hearing impairments.

- **Voice Assistants:** ML-powered voice assistants like Siri, Google Assistant, and Alexa utilize speech recognition to understand and respond to user queries, perform tasks, provide information, and control smart devices. ML enables natural language understanding, voice activation, and contextual responses.
- **Voice Search:** ML algorithms are employed to convert spoken search queries into text and retrieve relevant search results. Voice search applications are used in search engines, e-commerce platforms, and voice-enabled mobile applications.
- **Speech Analytics:** ML techniques are utilized to analyze large volumes of recorded speech data for customer service interactions, market research, sentiment analysis, and voice of the customer insights. ML algorithms can identify keywords, detect emotions, and extract valuable information from spoken conversations.
- **Automatic Speech Recognition (ASR):** ML models are trained to automatically transcribe spoken language into text. ASR systems find applications in call centers, voice-controlled systems, real-time captioning, and language learning platforms.
- **Speaker Identification and Verification:** ML algorithms can identify and verify the identity of a speaker based on their voice characteristics. These techniques are used in security systems, access control, and forensic investigations.
- **Language and Accent Adaptation:** ML models can adapt and recognize different languages and accents by learning from diverse speech data. ML algorithms enable multilingual speech recognition, accent conversion, and dialect adaptation.
- **Speech Emotion Recognition:** ML techniques can analyze speech signals to recognize and classify emotions conveyed in spoken language. This has applications in customer sentiment analysis, mental health assessment, and human-computer interaction.

These are some of the prominent applications of ML in speech recognition. ML algorithms such as Hidden Markov Models (HMMs), Recurrent Neural Networks (RNNs), Convolutional Neural Networks (CNNs), and Transformers are commonly used in building speech recognition systems.