

What is Reinforcement Learning?

AES ARP VQ

Reinforcement Learning (RL) is a type of machine learning where an agent learns how to behave in an environment by performing actions and receiving feedback in the form of rewards or penalties. The goal of the agent is to maximize the cumulative reward over time.

RL differs from other types of learning like supervised and unsupervised learning in that it focuses on learning by interaction rather than learning from a fixed dataset.

Reinforcement Learning (RL) is built upon several key elements that define how an agent interacts with its environment to learn optimal behaviors. These elements form the foundation of RL systems and help in formulating and solving problems. Let's dive into the **elements of reinforcement learning** in detail:

1. Agent

- **Definition:**

The **agent** is the decision-maker in RL. It interacts with the environment, observes its state, takes actions, and learns from the rewards it receives.

- **Role:**

The agent's primary objective is to find the best policy (strategy) that maximizes cumulative rewards over time.

Example:

In a self-driving car scenario, the car itself acts as the agent.

2. Environment

- **Definition:**

The **environment** represents everything the agent interacts with. It defines the rules and dynamics of the system, including how states change and rewards are provided.

- **Role:**

It provides feedback to the agent based on the actions taken.

Example:

In a chess game, the chessboard and its rules form the environment.

3. State (S)

- **Definition:**

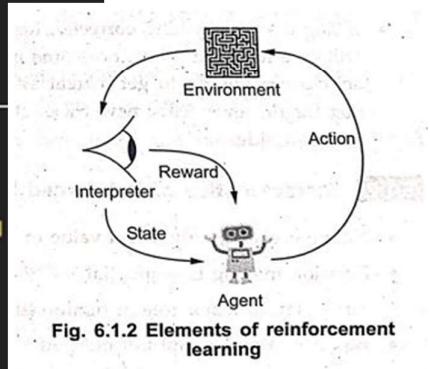
A **state** is a representation of the environment's condition at a particular time. It contains all the information the agent needs to make decisions.

- **Role:**

States help the agent understand "where" it is in the environment.

Example:

In a robot navigation problem, the robot's location and orientation in a maze define the state.



4. Action (A)

- **Definition:**

An **action** is a move or decision taken by the agent in a given state.

- **Role:**

The agent chooses actions to interact with the environment and progress toward its goal.

Example:

In a self-driving car, actions include accelerating, braking, or turning left/right.

5. Reward (R)

- **Definition:**

A **reward** is the feedback signal received by the agent after taking an action. It indicates how good or bad the action was in the given state.

- **Role:**

Rewards guide the agent's learning process by encouraging desirable actions and discouraging undesirable ones.

Example:

- In a game, winning gives a reward of +10, and losing gives a penalty of -10.
- A robot receives a reward for successfully reaching its goal without hitting obstacles.

6. Policy (π)

- **Definition:**

A **policy** is the agent's strategy or decision-making function. It defines the mapping from states to actions.

- **Role:**

The policy determines the agent's behavior. It can be deterministic (specific action for a state) or stochastic (probabilities for different actions).

Example:

- A policy in a chess-playing agent may dictate moving a pawn forward when in a specific board configuration.

7. Value Function (V)

- Definition:

The value function estimates the long-term reward of being in a particular state, assuming the agent follows a specific policy.

- Role:

It helps the agent evaluate how "good" a state is in terms of future rewards.

Formula:

$$V(s) = \mathbb{E}[R_t + R_{t+1} + R_{t+2} + \dots | s_t = s, \pi]$$

where R_t is the reward at time t .

Example:

In a maze, a state closer to the goal has a higher value than one farther away.

8. Q-Value (Action-Value Function, Q(S, A))

- Definition:

The Q-value (action-value function) estimates the long-term reward of taking a specific action in a given state and then following a particular policy.

- Role:

It helps the agent evaluate which actions are better in a given state.

Formula:

$$Q(s, a) = \mathbb{E}[R_t + \gamma V(s_{t+1}) | s_t = s, a_t = a]$$

where γ is the discount factor.

Example:

In a game, the Q-value might evaluate the benefits of making a risky move versus a safe one.

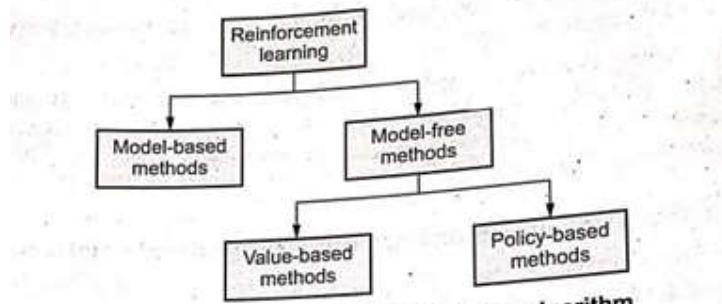


Fig. 6.2.1 Types of Reinforced learning algorithm

1. Model-Based Methods

- **Definition:**
 - These methods involve building a model of the environment. The model predicts the behavior of the environment, including the next state and the reward, given a current state and an action.
- **Key Characteristics:**
 - Used for planning by simulating future steps.
 - Focuses on utilizing or improving the model to derive the optimal policy.
- **Example Algorithms:**
 - Dynamic Programming (Value Iteration, Policy Iteration).
 - Monte Carlo Tree Search (MCTS).

2. Model-Free Methods

- **Definition:**
 - These methods do not require the agent to learn or use a model of the environment. The agent learns directly from its interactions with the environment.
- **Subdivided into:**
 1. **Value-Based Methods:**
 - The agent focuses on learning value functions like $V(s)$ (state value) or $Q(s, a)$ (action value).
 - **Key Idea:**
 - The policy is indirectly derived from the value function by selecting actions with the highest expected reward.
 - **Example Algorithms:**
 - Q-Learning.
 - SARSA (State-Action-Reward-State-Action).

2. Policy-Based Methods:

- The agent directly learns a policy $\pi(a|s)$, which maps states to actions, without estimating value functions.
- **Key Idea:**
 - The policy is optimized to maximize the expected cumulative reward using techniques like Policy Gradients.
- **Example Algorithms:**
 - REINFORCE algorithm.
 - Policy Gradient Methods.

Summary of the Flow

- Reinforcement Learning is broadly categorized into:
 1. **Model-Based Methods:** Uses a model of the environment to plan.
 2. **Model-Free Methods:** Relies purely on experience without modeling.
 - **Value-Based:** Focuses on value functions to guide actions.
 - **Policy-Based:** Directly optimizes the policy for decision-making.

1. Positive vs. Negative Reinforcement:

- **Positive Reinforcement:** Strengthens behavior by providing rewards.
- **Negative Reinforcement:** Strengthens behavior by removing negative outcomes.

Real-Time Applications of Reinforcement Learning

1. **Robotics:**
 - **Application:** RL helps robots learn tasks like walking, picking objects, or performing assembly operations.
 - **Example:** Training robotic arms to adapt to new tasks in manufacturing.
2. **Autonomous Vehicles:**
 - **Application:** RL is used for decision-making and control in self-driving cars.
 - **Example:** Optimizing navigation, lane-keeping, and obstacle avoidance.
3. **Game Playing:**
 - **Application:** Training agents to play games at superhuman levels.
 - **Example:** AlphaGo defeating human champions in Go.
4. **Healthcare:**
 - **Application:** Optimizing treatment plans, drug discovery, and patient care strategies.
 - **Example:** Personalized medicine by finding the best sequences of treatments.
5. **Finance:**
 - **Application:** Portfolio optimization, trading strategies, and fraud detection.
 - **Example:** RL-driven stock trading systems adapting to market changes.
6. **Manufacturing:**
 - **Application:** Managing supply chains, optimizing production lines, and reducing energy consumption.
 - **Example:** Predictive maintenance scheduling.
7. **Energy Systems:**
 - **Application:** Balancing supply and demand in energy grids.
 - **Example:** RL in smart grid management to minimize energy wastage.
8. **Marketing:**
 - **Application:** Personalizing user recommendations, ad placements, and customer retention strategies.
 - **Example:** Optimizing ad display for higher click-through rates.
9. **Natural Language Processing:**
 - **Application:** Chatbot training, conversational AI, and machine translation.
 - **Example:** RL-based dialogue systems that improve user satisfaction.
10. **Traffic Management:**
 - **Application:** Optimizing traffic signals and reducing congestion.
 - **Example:** Smart traffic light systems using RL to adapt to traffic flow.

Need for Reinforcement Learning (RL)

Reinforcement Learning (RL) is a critical approach in machine learning for solving problems where decision-making is required in dynamic, uncertain, and complex environments. Below are detailed reasons explaining why RL is needed:

2. Learning from Interaction (Trial-and-Error)

- Unlike supervised learning, where labeled data is required, RL allows agents to learn directly from interaction with the environment using feedback (rewards or penalties).
- Example:
 - A robot learning to walk: Instead of being programmed explicitly, it learns by trying movements, falling, and adjusting actions to improve.

Why RL?

RL is needed when explicit programming of rules is infeasible or when data labeling is impractical or impossible. It enables systems to learn autonomously.

3. Handling Unknown or Complex Environments

- Many real-world environments are too complex to model accurately or are entirely unknown beforehand.
- Example:
 - In robotics, modeling every detail of an environment (e.g., obstacles, friction) is impractical.

Why RL?

RL allows agents to operate effectively without a detailed model of the environment, making decisions based on observed feedback rather than relying on perfect information.

4. Dynamic Environments

- RL is designed to handle environments that change over time. Traditional optimization or rule-based systems struggle in such scenarios.
- Example:
 - Stock market trading strategies need to adapt to changing market conditions.

Why RL?

RL agents can continually learn and adapt policies based on new data, making them suitable for evolving environments.

1. **Sequential Decision-Making :** Reinforcement learning is designed to tackle problems that involve sequential decision-making, where actions have consequences and outcomes unfold over time. In such scenarios, traditional machine learning approaches that rely on static datasets or supervised learning may not be suitable. Reinforcement learning allows agents to learn and adapt their behaviour based on feedback received from the environment on random as-you-go basis.

7. Automation of Complex Tasks

- Many complex tasks, such as controlling drones or industrial robots, cannot be fully automated using predefined rules or supervised learning.
- Example:
 - In manufacturing, robots may need to adapt to different tasks, such as assembling parts of varying shapes.

Why RL?

RL enables the automation of such tasks by learning optimal behaviors through interaction, reducing the need for manual intervention or rule-setting.

10. Applications Requiring Real-Time Learning

- In some applications, decisions need to be made in real-time, with continuous learning and adaptation.
- Example:
 - Autonomous cars need to learn and make decisions in real-time based on traffic conditions, road layouts, and obstacles.

Why RL?

RL provides a robust framework for learning and decision-making in real-time scenarios.

2. **Limited or Unavailable Labelled Data :** In many real-world scenarios, obtaining labelled data for training machine learning models can be challenging, time-consuming, or even impossible. Reinforcement learning provides a framework where agents can learn from interactions with the environment, acquiring knowledge and improving performance through trial and error, without relying on explicitly labelled data.

5. Exploration vs. Exploitation Trade-off

- In many scenarios, systems must decide whether to exploit known rewarding actions or explore new actions that might yield better rewards.
- Example:
 - In a recommendation system, should a new movie be recommended to explore user preferences, or should a proven favorite movie be recommended to ensure satisfaction?

Why RL?

RL provides mechanisms (e.g., ϵ -greedy policies) to balance exploration and exploitation, which is essential for long-term optimization.

Aspect	Supervised Learning	Unsupervised Learning	Reinforcement Learning
Definition	Learning from labeled data to make predictions or classify data.	Learning patterns and structures from unlabeled data.	Learning by interacting with an environment and maximizing cumulative rewards.
Goal	Predict the output (labels) for given inputs.	Discover hidden structures, patterns, or relationships in data.	Learn an optimal policy for decision-making over time.
Input Data	Labeled data (input-output pairs).	Unlabeled data (no predefined labels).	Environment states and actions; no explicit input-output pairs.
Output	Predictive models for regression or classification tasks.	Clusters, patterns, or dimensionality-reduced representations of data.	A policy that maps states to actions for maximum reward.
Feedback	Explicit (labels guide learning).	No explicit feedback during training.	Feedback through rewards or penalties based on actions.
Learning Process	Directly maps inputs to outputs using labels.	Finds structure or similarity in data (e.g., clusters, features).	Learns from trial-and-error through interaction with the environment.
Algorithms	Linear Regression, Logistic Regression, Decision Trees, Neural Networks.	K-Means, DBSCAN, PCA, Autoencoders, Hierarchical Clustering.	Q-Learning, Deep Q-Learning (DQN), Policy Gradient, Actor-Critic Methods.
Key Applications	Spam detection, image recognition, stock price prediction, speech recognition.	Market segmentation, anomaly detection, recommendation systems.	Robotics, gaming (AlphaGo, chess), autonomous vehicles, stock trading.
Strengths	Accurate predictions if sufficient labeled data is available.	Useful for exploring data without labels, reducing dimensionality, and anomaly detection.	Adapts to dynamic environments and learns optimal strategies over time.
Challenges	Requires large labeled datasets; time-consuming labeling process.	Lack of interpretability; sensitive to noise and outliers.	High computational cost; balancing exploration and exploitation.
Nature of the Problem	Static: Based on fixed datasets.	Static: Analyzes data without considering sequential actions.	Dynamic: Focuses on sequential decision-making in changing environments.

i) Supervised Learning

Definition:

Supervised Learning is a type of machine learning where the model learns from labeled data. The goal is to map input data to the corresponding output (target) by identifying patterns in the data.

Key Characteristics:

- **Labeled Data:** Each input in the training dataset has a corresponding output label.
- **Objective:** Predict the output for new, unseen inputs.
- **Types:**
 - **Regression:** Predict continuous values (e.g., stock prices).
 - **Classification:** Predict discrete categories (e.g., spam detection).

Examples:

1. Predicting house prices based on features like area, number of rooms, etc.
2. Classifying emails as spam or not spam.

Common Algorithms:

- Linear Regression, Logistic Regression.
- Decision Trees, Random Forest.
- Support Vector Machines (SVM).
- Neural Networks.

ii) Unsupervised Learning

Definition:

Unsupervised Learning is a type of machine learning where the model works with unlabeled data to discover patterns, groupings, or structures in the dataset.

Key Characteristics:

- **Unlabeled Data:** No predefined outputs or labels.
- **Objective:** Identify hidden patterns, group similar data points, or reduce dimensionality.
- **Types:**
 - **Clustering:** Grouping similar data points (e.g., customer segmentation).
 - **Dimensionality Reduction:** Reducing the number of features while preserving important information (e.g., Principal Component Analysis).

Examples:

1. Market segmentation: Grouping customers based on purchasing behavior.
2. Anomaly detection: Identifying fraudulent transactions in financial data.

Common Algorithms:

- K-Means Clustering, DBSCAN.
- Hierarchical Clustering.
- PCA (Principal Component Analysis).

iii) Reinforcement Learning

Definition:

Reinforcement Learning (RL) is a type of machine learning where an agent interacts with an environment, learns from feedback (rewards or penalties), and aims to maximize cumulative rewards over time.

Key Characteristics:

- **Agent and Environment:** The agent takes actions in an environment to achieve a goal.
- **Trial-and-Error:** The agent learns through experimentation, receiving rewards or penalties based on actions.
- **Sequential Decision-Making:** RL focuses on optimizing actions in a sequence, considering future consequences.

Examples:

1. Training a robot to walk by rewarding stable movements.
2. Playing games like chess or Go, where the agent learns strategies to win.
3. Optimizing stock trading strategies by maximizing profits over time.

Common Algorithms:

- Q-Learning, Deep Q-Learning (DQN).
- Policy Gradient Methods (e.g., PPO, REINFORCE).
- Actor-Critic Methods (e.g., A3C).

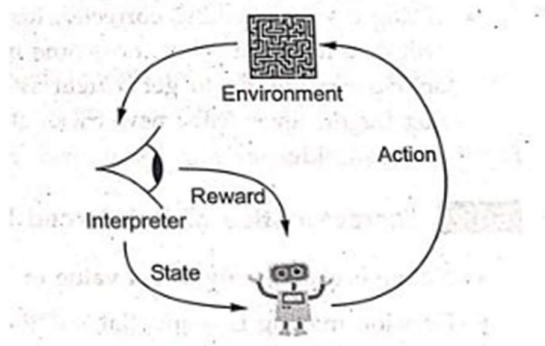


Fig. 6.1.2 Elements of reinforcement learning

The Markov Property

It defines the relationship between a state and its future, ensuring that the future state depends only on the current state and not on the sequence of states that preceded it. This is often referred to as the "memoryless" property because the past is irrelevant to predicting the future, given the present.

Formal Definition:

The Markov Property states that the probability of transitioning to the next state depends only on the current state and action, not on the sequence of events that led to the current state. Formally, it can be written as:

$$P(s_{t+1}|s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_1, a_1) = P(s_{t+1}|s_t, a_t)$$

Here:

- s_t is the state at time t ,
- a_t is the action taken at time t ,
- s_{t+1} is the state at the next time step.

This equation expresses that the future state s_{t+1} depends only on the current state s_t and the current action a_t , and not on past states or actions.

Intuition Behind the Markov Property:

- **Memoryless:** The agent only needs to know the present state to make predictions about future states. It doesn't need to store or recall information about earlier states or actions.
- **Simplifies Modeling:** This property significantly reduces the complexity of modeling a system, as you don't need to consider the entire history of the agent's interaction with the environment, only the present state.
- **State Encapsulation:** A state encapsulates all relevant information about the system needed for future predictions, meaning that the environment's dynamics can be described entirely by the current state.

Markov Assumptions

Markov processes (and models based on these processes) make the following assumptions.

1. The number of possible outcomes or states is finite.
2. The outcome at any stage depends only on the outcome of the previous stage.
3. The probabilities of transitioning from one state to another state are constant over time.

Example:

Consider a simple game where an agent moves on a grid:

- The agent's state is the current position on the grid.
- The action is to move in one of four directions: up, down, left, or right.

With the Markov property, the agent's future position (state) depends only on its current position and the action it takes, not on how it arrived at that position. For example, if the agent is at position (2, 3) and decides to move right, the next state (position) will be (2, 4), and the previous moves or positions are irrelevant to determining that outcome.

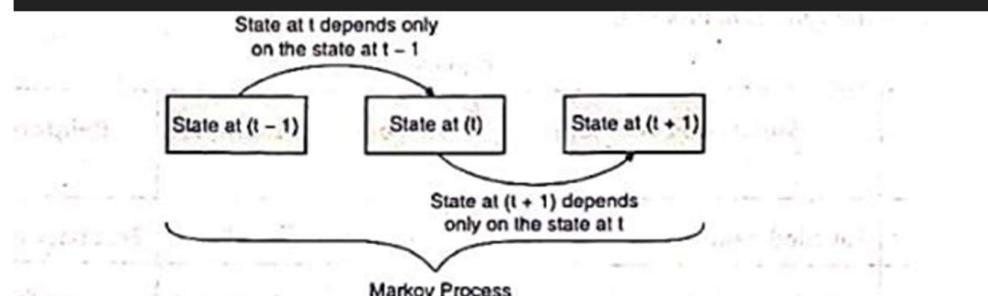


Fig. 6.2.1 : Markov Process

Markov Chain

Definition : A Markov chain is a mathematical model that represents the state transitioning probabilities of a Markov process.

- the next state of a Markov process depends only upon the current state. Markov chain assigns state transitioning probabilities for a Markov process.

3. Stochastic Process: A Markov chain is a type of random process where the state transitions are probabilistic. This means that at each time step, the chain moves to one of the possible next states with certain probabilities.

Transition Matrix:

A Markov chain's behavior can be described using a **transition matrix** P , where each element P_{ij} represents the probability of transitioning from state i to state j .

For example, consider a system with three states S_1, S_2, S_3 . The transition matrix might look like this:

$$P = \begin{bmatrix} P_{11} & P_{12} & P_{13} \\ P_{21} & P_{22} & P_{23} \\ P_{31} & P_{32} & P_{33} \end{bmatrix}$$

Each row sums to 1, as the agent must transition to some state after each time step.

- For example, the Fig. 6.2.2 illustrates a simple Markov chain for weather.

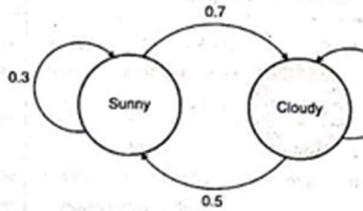


Fig. 6.2.2 : Markov chain for weather

A **Markov Reward Process (MRP)** is a type of Markov process that extends the standard **Markov Chain** by incorporating **rewards** into the model. In an MRP, each state transition is associated with a **reward**, and the process is used to evaluate how an agent performs in a given environment where states and transitions have associated rewards.

Goal of MRP:

The goal in an MRP is often to compute the **value function** for each state, which represents the expected cumulative reward starting from that state, considering the discount factor.

Steps to Solve MRP:

1. **Initialization:** Start with an initial guess for the value function $V(s)$ for all states, often setting all values to zero.
2. **Iteration:** Update the value function for each state using the Bellman equation. This process can be done iteratively until the value function converges to a stable solution.
3. **Convergence:** The algorithm continues updating $V(s)$ until the changes are small enough (i.e., convergence criterion is met).

A **Markov Decision Process (MDP)** is an extension of a **Markov Reward Process (MRP)** that incorporates **actions** taken by an agent in an environment. It provides a formal framework for modeling decision-making problems where an agent interacts with an environment, making decisions to maximize a cumulative reward over time.

In an MDP, the agent can choose actions that influence the transition between states and the associated rewards, making it suitable for solving reinforcement learning problems, where the goal is for the agent to learn an optimal policy to maximize long-term rewards.

- So, the transition model of a MDP can be given as $p(s', r | s, a)$ which is read as the probability that the agent leaves the environment in state s' (after carrying out the action) and gets a reward r given the previous state s and action a .
- You can think of MDP as an interface. It is a contract between the agent and the environment. The agent can only observe the states and suggest an action. In return, the environment gives the agent a new state and some value that represents how well the agent is doing. The environment and the agent can be as complex as they need to be.

The Bellman equation is a key recursive relationship used in Markov Decision Processes (MDPs) to compute the value function of states and actions, which is central to solving reinforcement learning problems. The equation expresses the relationship between the value of a state (or state-action pair) and the expected rewards that can be obtained from that state, given an agent's policy.

The Bellman equation helps break down a decision-making process into smaller subproblems, which makes it possible to compute the optimal policy for an agent.

- In simple terms, the Bellman equation states that the value of being in a particular state (or taking a specific action in a state) is equal to the immediate reward obtained from that state (or state-action pair) plus the expected value of being in the next state (or the expected value of taking the next action in the next state).

- Mathematically, the Bellman equation can be written as follows:

$$V(s) = R(s) + \gamma \sum P(s, a, s') * V(s')$$

In this equation:

- $V(s)$ represents the value of state s , which indicates how desirable or valuable it is to be in that state.
 - $R(s)$ is the immediate reward obtained when transitioning to state s .
 - γ (gamma) is a discount factor that determines the importance of future rewards compared to immediate rewards. It is a value between 0 and 1.
 - $P(s, a, s')$ is the probability of transitioning from state s to state s' when taking action a .
- 💡 Essentially, the Bellman equation tells you that the value of a state is the sum of the immediate reward and the discounted expected value of the next state. By iteratively applying the Bellman equation to update the value estimates of states or state-action pairs, you can eventually converge to the optimal values that maximise the cumulative rewards.
- The Bellman equation is a crucial tool for solving reinforcement learning problems and finding optimal policies. It allows you to break down a complex decision-making process into simpler steps and determine the value of each state or state-action pair based on future rewards and transitions.

Example:

In a simple grid world, where an agent can move around to reach a goal, the Bellman equation is used to calculate the value of each state considering possible actions. If the agent is at a particular location, the Bellman equation allows it to calculate the expected reward from that state, considering all possible actions (e.g., up, down, left, right) and the resulting transitions to new states, until the goal is reached.

Q-Learning is a model-free reinforcement learning algorithm used to learn the optimal action-selection policy in a given environment. It enables an agent to maximize the total reward over time by learning which actions to take in various states, without requiring prior knowledge of the environment.

It does this by iteratively updating a table of Q-values, where each Q-value represents the expected future reward for taking a specific action in a given state and following the optimal policy thereafter.

1. Agent

The entity that interacts with the environment, learns from it, and makes decisions.

- Example: A robot navigating a maze.

2. Environment

The external system with which the agent interacts. It provides feedback (rewards) based on the agent's actions.

- Example: The maze or game world.

3. State (S)

A representation of the current situation or position of the agent in the environment.

- Example: A cell in a grid (e.g., (2, 3)).

4. Action (A)

The choices available to the agent in a given state. Actions determine how the agent interacts with the environment.

- Example: Move up, down, left, or right.

5. Reward (R)

The feedback received from the environment after performing an action. Rewards guide the agent to learn the best behavior.

- Example:
 - +10: For reaching the goal.
 - -1: For hitting a wall.

6. Q-Value ($Q(s, a)$)

The expected future reward for taking action a in state s and following the optimal policy thereafter. It reflects the "quality" of a state-action pair.

7. Policy (π)

A strategy that defines which action to take in a given state. Q-Learning indirectly learns the optimal policy by improving Q-values over time.

8. Learning Rate (α)

A hyperparameter that controls how much new information overrides the old information when updating Q-values.

- Range: $0 < \alpha \leq 1$.

9. Discount Factor (γ)

A hyperparameter that determines the importance of future rewards compared to immediate rewards.

- Range: $0 \leq \gamma < 1$.
- Higher γ : Focus on long-term rewards.
- Lower γ : Focus on short-term rewards.

11. Q-Table

A table used to store Q-values for each state-action pair.

- Rows: States (S).
- Columns: Actions (A).

12. Bellman Equation

The recursive formula used to update Q-values:

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

It combines immediate rewards (R) with the estimated future rewards ($\max_{a'} Q(s', a')$).

The **Q-Table** is a core component of Q-Learning. It is a **tabular representation** that stores the **Q-values** ($Q(s, a)$) for all possible state-action pairs in an environment. These Q-values represent the agent's estimate of the expected future rewards for taking a specific action (a) in a particular state (s) and following the optimal policy thereafter.

Structure of the Q-Table

- **Rows:** Represent all possible states (S) in the environment.
- **Columns:** Represent all possible actions (A) the agent can take.
- **Cells:** Contain the Q-values ($Q(s, a)$), which indicate the "quality" of taking action a in state s .

Purpose of the Q-Table

The Q-Table helps the agent:

1. Learn and store the best actions to take in each state over time.
2. Make decisions by choosing the action with the highest Q-value for the current state ($\max_a Q(s, a)$).
3. Gradually converge toward the optimal policy as it updates its estimates.

Initialization of the Q-Table

Initially:

- The Q-Table is filled with **zeros** (or small random values) because the agent has no prior knowledge of the environment.
- The table is updated as the agent explores the environment and learns through rewards and penalties.

Updating the Q-Table

The Q-values in the table are updated using the **Q-Learning update rule** (Bellman equation):

The **Q-function** is a key concept in reinforcement learning, especially in Q-Learning. It is a **mathematical function** that quantifies the expected cumulative reward of an agent starting from a given state s , taking a specific action a , and following the optimal policy thereafter.

The Q-function is often denoted as:

$$Q(s, a)$$

- $Q(s, a)$ represents the **quality** of taking action a in state s .
- Specifically, it is the **expected total reward** the agent will receive over time, starting from state s , taking action a , and then acting optimally in subsequent states.

Purpose of the Q-Function

The Q-function helps the agent:

1. **Evaluate Actions:** Determines which action a is better in a given state s .
2. **Optimal Policy Derivation:** The optimal action in any state is the one with the **highest Q-value**:

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

This leads to the optimal policy, π^* .

Core Idea of Q-Learning

The agent interacts with an environment, receives feedback in the form of rewards, and learns to make better decisions over time. The **Q-values** stored in a **Q-Table** are iteratively updated based on the feedback, eventually converging to the optimal policy.

Steps in Q-Learning Algorithm

1. Initialization

- Initialize the **Q-Table**:
 - The Q-Table is a matrix where rows represent states (s) and columns represent actions (a).
 - Each entry $Q(s, a)$ starts with a default value (e.g., 0).
- Define the hyperparameters:
 - **Learning rate (α)**: Controls how much new information overrides old information ($0 < \alpha \leq 1$).
 - **Discount factor (γ)**: Balances the importance of immediate vs. future rewards ($0 \leq \gamma < 1$).
 - **Exploration probability (ϵ)**: Balances exploration (trying new actions) and exploitation (choosing the best-known action).

2. Repeat for Each Episode

An episode is a complete interaction with the environment, ending in a terminal state or after a set number of steps.

1. Start at an Initial State (s):

- The agent begins in a specific state of the environment.

2. Choose an Action (a):

- Use the **ϵ -greedy policy** to select an action:
 - **Exploration**: With probability ϵ , choose a random action to explore the environment.
 - **Exploitation**: With probability $1 - \epsilon$, choose the action with the highest Q-value ($\arg \max_a Q(s, a)$).

3. Perform the Action:

- Execute the selected action a and observe:
 - The reward (R) provided by the environment.
 - The next state (s') the agent transitions to.

4. Update the Q-Value:

- Use the Q-Learning update rule (Bellman equation):

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

- $Q(s, a)$: Current Q-value.
- α : Learning rate.
- R : Reward for taking action a in state s .
- γ : Discount factor.
- $\max_{a'} Q(s', a')$: Maximum Q-value for the next state s' , representing the best future action.

5. Transition to the Next State (s'):

- Update the agent's state to s' .

6. Repeat Until the Episode Ends:

- Continue selecting actions, updating Q-values, and transitioning between states until the episode ends.

3. Convergence

- After many episodes, the Q-values converge to the optimal Q-function ($Q^*(s, a)$), enabling the agent to derive the optimal policy:

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

Example: Q-Learning in a Grid World

Environment:

- **Grid:** A 4x4 grid with the agent starting at the top-left cell (0,0) and the goal at the bottom-right cell (3,3).
- **Actions:** Move up, down, left, right.
- **Rewards:**
 - +10 for reaching the goal.
 - -1 for hitting a wall.
 - 0 for other transitions.

Algorithm Execution:

1. Initialize Q-Table:
 - All Q-values are set to 0 initially.
2. Agent explores the grid:
 - Takes actions, receives rewards, and updates Q-values using the Q-Learning formula.
3. Over multiple episodes:
 - The agent learns the optimal path to the goal (e.g., shortest path).