# Reinforcement and Deep Learning

## Unit-1

---

### Introduction to Reinforcement Learning (RL)

**Reinforcement Learning (RL)** is a subfield of machine learning focused on how agents should take actions in an environment to maximize cumulative rewards over time. Unlike supervised learning, where the model learns from labeled data, RL is concerned with learning from interaction. The core idea is learning through **trial and error**, where an agent repeatedly takes actions, observes outcomes, and adjusts its strategy to improve its performance.

### 1. The RL Problem

At its core, RL deals with decision-making in uncertain environments. The agent interacts with the environment by:

- **Perceiving the state** of the environment.
- **Taking an action** based on the current state.
- **Receiving feedback (reward)** from the environment based on the action taken.
- **Updating its strategy** (policy) to improve future actions and maximize the cumulative reward.

The process continues iteratively, and the agent's goal is to **maximize the total reward** (also known as the return) over time, typically discounted over future actions to prioritize immediate rewards.

### 2. Reinforcement Learning Cycle

The basic RL cycle follows a simple feedback loop:

1. **The agent observes the current state (s)** of the environment.
2. **The agent selects an action (a)** according to its policy (π).
3. **The environment transitions to a new state (s')** based on the action.
4. **The environment gives a reward (r)** to the agent.
5. **The agent updates its policy** based on the reward and transitions to the next state.

This cycle continues, and the agent uses these interactions to gradually improve its decision-making strategy (policy).

### 3. Exploration vs. Exploitation

One of the key challenges in RL is the trade-off between:

- **Exploration**: Trying new actions to discover their rewards, which may lead to better long-term outcomes. However, exploring unknown actions may result in lower immediate rewards.
- **Exploitation**: Using the current best-known actions to maximize rewards in the short term. However, this may prevent the agent from discovering potentially better actions.

Balancing these two is crucial, as too much exploration can lead to poor performance, while too much exploitation may result in the agent missing out on better long-term rewards.

---

**Key Terms in RL**:

- **Agent**: The learner or decision-maker that interacts with the environment. The agent chooses actions and seeks to maximize rewards.
- **Environment**: The system that the agent interacts with. The environment provides feedback in the form of rewards or penalties after each action taken by the agent.
- **State (s)**: A snapshot of the current situation in the environment. It represents everything that the agent needs to know to make a decision. States can be simple (e.g., a position on a grid) or complex (e.g., pixel values in an image).
- **Action (a)**: A choice made by the agent that affects the state of the environment. For example, in a game, an action could be moving left, right, jumping, etc.
- **Reward (r)**: The immediate feedback given to the agent by the environment after taking an action. Positive rewards encourage the behavior, while negative rewards discourage it.
- **Policy ($\pi$)**: The agent's strategy for deciding which actions to take in different states. It maps states to actions and can be either deterministic (always taking the same action in a given state) or stochastic (choosing actions based on probabilities).
- **Return (G)**: The cumulative reward that the agent seeks to maximize. This may include future rewards, often discounted by a factor to prioritize immediate gains.
- **Value Function (V)**: The **value function** estimates the expected cumulative reward (return) for each state. It indicates how good it is for the agent to be in a particular state under a certain policy.

---

**Features of RL**

- **Learning from Interaction**: In RL, an agent learns by interacting with its environment. Unlike supervised learning, where the agent learns from labeled data, RL relies on trial and error, adjusting its behavior based on the outcomes of its actions.

- **Delayed Rewards**: In RL, rewards are not always immediate. The agent may need to perform a sequence of actions before receiving a reward. The goal is to maximize the long-term reward rather than the immediate gain, which involves handling **delayed gratification**.

- **Exploration and Exploitation Trade-off**: RL involves balancing **exploration** (trying new actions to discover their effects) and **exploitation** (using the current knowledge to maximize immediate rewards). Too much exploration can waste time on ineffective actions, while too much exploitation can miss out on potentially better strategies.

- **Trial and Error Learning**: RL is inherently a process of trial and error. The agent explores the environment, tries different actions, and learns from the outcomes

(rewards or penalties). Over time, the agent improves its decisions to maximize long-term rewards.

- **Sequential Decision-Making**:The agent makes decisions sequentially over time. Each action affects future actions and outcomes, so RL is about planning and making decisions that optimize long-term outcomes, not just immediate rewards.

- **Model-Free and Model-Based Methods**: **Model-free** methods (like Q-learning) learn directly from experience without building a model of the environment.

  **Model-based** methods (like Dynamic Programming) involve building a model of the environment to predict future states and rewards.

- **Markov Property**: RL systems are often modeled as **Markov Decision Processes (MDPs)**, which assume that the future depends only on the current state and action (not the history of states). This is called the **Markov property**. It simplifies decision-making by focusing only on the present situation.

---

# Reinforcement Learning Framework and Markov Decision Process (MDP)

## Markov Decision Process (MDP)

A **Markov Decision Process (MDP)** is a mathematical framework used to model decision-making problems where outcomes are partly random and partly under the control of a decision-maker (agent). MDPs are commonly used to formalize RL problems because they describe the environment in which an agent operates.

**Key Components of MDP**:

1. States (S):

   - **States** represent all possible situations or configurations the environment can be in. At any given time, the environment is in a specific state $s$. The agent observes this state to decide which action to take.

   - Example: In a chess game, a state could represent the current arrangement of pieces on the board.

2. Actions (A):

   - **Actions** are the choices or moves available to the agent. For each state, the agent selects an action to influence the environment.

   - The set of actions available can be discrete (e.g., up, down, left, right) or continuous (e.g., adjusting the throttle of a car).

3. **Transition Function (P or T):**

   - The **transition function** defines the probability of moving from one state to another after taking a certain action. It is denoted as $P(s'|s, a)$, where $s'$ is the next state, $s$ is the current state, and $a$ is the action taken.

   - The transition function models the dynamics of the environment.

   - **Markov Property:** The transition to the next state depends only on the current state and action, not on the history of previous states. This simplifies the decision-making process because only the present matters.

4. **Reward Function (R):**

   - The **reward function** defines the immediate reward received by the agent after transitioning from state $s$ to state $s'$ via action $a$. It is denoted as $R(s, a, s')$.

   - The reward function provides feedback on the quality of the agent's actions and guides the agent toward maximizing the cumulative reward.

5. **Policy (π):**

   - A **policy (π)** is the strategy that defines the agent's behavior. It is a mapping from states to actions. The policy can be:

     - **Deterministic:** Always selects the same action for a given state.

     - **Stochastic:** Selects actions based on a probability distribution over actions.

   - The goal is to find the **optimal policy (π\*)** that maximizes the expected cumulative reward.

6. **Value Function (V(s)) and Action-Value Function (Q(s, a)):**

   - The **value function** estimates the expected return from a given state. The value function for a policy $\pi$ is denoted as:

$$V^{\pi}(s) = \mathbb{E}[G_t|s_t = s, \pi]$$

   where $G_t$ is the cumulative reward (return) from time step $t$ onward.

   - The **action-value function** $Q(s, a)$ estimates the expected return for taking action $a$ in state $s$, and then following policy $\pi$:

$$Q^{\pi}(s, a) = \mathbb{E}[G_t|s_t = s, a_t = a, \pi]$$

   - These functions are used to evaluate and improve the policy over time.

7. **Discount Factor (γ):**

- The **discount factor** $\gamma$ is a value between 0 and 1 that determines the importance of future rewards. A value close to 1 means future rewards are nearly as important as immediate rewards, while a value close to 0 means the agent prioritizes immediate rewards.

**Markov Property:**

- The **Markov property** states that the future state depends only on the current state and the action taken, not on the sequence of states that preceded it. This simplifies decision-making by ensuring the agent only needs to consider the present state when making a choice.

**Goal of MDP:**

- The primary goal of an MDP is to find an **optimal policy** $\pi\backslash^*$ that maximizes the expected cumulative reward (or return) over time. The agent achieves this by interacting with the environment, learning from feedback, and adjusting its policy based on experience.

$\downarrow$

## Solving MDPs:

To solve an MDP and find the optimal policy, different algorithms can be used:

1. **Dynamic Programming**: Uses a model of the environment (transition and reward functions) to compute the optimal policy.

2. **Monte Carlo Methods**: Learn the value of states and actions through sample episodes without requiring a model of the environment.

3. **Temporal-Difference Learning (TD)**: Combines ideas from Dynamic Programming and Monte Carlo methods by updating value estimates based on the difference between consecutive estimates.

---

# Policies, Value Functions, and Bellman Equations

## 1. Policies (π)

In reinforcement learning, a **policy** defines the agent's strategy for selecting actions at any given state. It serves as the decision-making rule, determining which action the agent will take based on the current state.

**Types of Policies:**

1. **Deterministic Policy:**

   - A **deterministic policy** maps each state directly to a single action. For any given state $s$, the agent always selects the same action $a$ under a deterministic policy.

   $$\pi(s) = a$$

   This means that for each state, the action is fixed.

2. **Stochastic Policy:**

   - A **stochastic policy** defines a probability distribution over actions for each state. The agent selects actions based on probabilities rather than deterministically.

   $$\pi(a|s) = P(a|s)$$

   This means that the action taken is random, but with probabilities assigned by the policy. For example, in state $s$, the agent might have a 70% chance of taking action $a_1$ and a 30% chance of taking action $a_2$.

## 2. Value Functions

**Value functions** are critical for evaluating how good it is for an agent to be in a particular state or to take a specific action. The value functions estimate the expected return (cumulative reward) from states and actions, helping the agent make informed decisions.

**Types of Value Functions:**

1. **State Value Function (V(s)):**

   - The **state value function** $V^{\pi}(s)$ measures the expected return (cumulative reward) the agent will receive starting from state $s$ and following policy $\pi$ thereafter.

   $$V^{\pi}(s) = \mathbb{E}_{\pi}[G_t|s_t = s]$$

   Here, $G_t$ is the return (sum of future discounted rewards), and $\mathbb{E}_{\pi}$ represents the expected value based on policy $\pi$.

   - In simpler terms, $V^{\pi}(s)$ tells us how good it is to be in a particular state $s$ when following policy $\pi$.

2. **Action Value Function (Q(s, a)):**

   - The **action-value function** $Q^{\pi}(s, a)$ measures the expected return starting from state $s$, taking action $a$, and then following policy $\pi$.

   $$Q^{\pi}(s, a) = \mathbb{E}_{\pi}[G_t|s_t = s, a_t = a]$$

- $Q^\pi(s, a)$ tells us how good it is to take a particular action $a$ in state $s$, and then follow the policy $\pi$ in subsequent steps.
- The **state-value function** $V^\pi(s)$ can be derived from the **action-value function** $Q^\pi(s, a)$:

$$V^\pi(s) = \sum_a \pi(a|s) Q^\pi(s, a)$$

This means that the value of a state is the expected value of the action-value function over all possible actions in that state, weighted by the probability of taking each action.

**Optimal Value Functions:**

- *Optimal State Value Function (V(s))**:
  - The **optimal state value function** represents the maximum expected return an agent can achieve starting from state $s$ and following the optimal policy $\pi^*$.

$$V^*(s) = \max_\pi V^\pi(s)$$

  - This tells us the best possible long-term reward that can be obtained from state $s$.
- *Optimal Action Value Function (Q(s, a))**:
  - The **optimal action value function** represents the maximum expected return starting from state $s$, taking action $a$, and then following the optimal policy $\pi^*$.

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

  - This tells us the best possible long-term reward obtainable by taking action $a$ in state $s$, and then acting optimally thereafter.

# 3. Bellman Equations

The **Bellman equations** provide recursive definitions for value functions. They are fundamental to many reinforcement learning algorithms, including Q-learning and Dynamic Programming. These equations describe the relationship between the value of a state (or action) and the values of its successor states (or actions).

**Bellman Equation for State Value Function:**

The **Bellman equation** for the state value function $V^\pi(s)$ expresses the value of a state as the immediate reward plus the discounted value of the next state, assuming the agent follows policy $\pi$.

$$V^\pi(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s, a)[R(s, a, s') + \gamma V^\pi(s')]$$

Where:

- $P(s'|s, a)$: The probability of transitioning from state $s$ to state $s'$ after taking action $a$.

- $R(s, a, s')$: The immediate reward received after transitioning from $s$ to $s'$ by taking action $a$.

- $\gamma$: The **discount factor**, determining the importance of future rewards.

This recursive equation says that the value of a state is the expected immediate reward plus the expected discounted value of the next state.

**Bellman Equation for Action Value Function**:

The **Bellman equation** for the action-value function $Q^\pi(s, a)$ similarly expresses the value of taking action $a$ in state $s$ as the immediate reward plus the discounted value of the next state, following policy $\pi$:

$$Q^\pi(s, a) = \sum_{s'} P(s'|s, a)[R(s, a, s') + \gamma \sum_{a'} \pi(a'|s')Q^\pi(s', a')]$$

This equation represents the expected return for taking action $a$ in state $s$ and then following policy $\pi$.

**Bellman Optimality Equations**:

For **optimal policies**, we use the **Bellman Optimality Equations**, which represent the best possible long-term return from any state or state-action pair:

- **State Value Function**:

$$V^*(s) = \max_a \sum_{s'} P(s'|s, a)[R(s, a, s') + \gamma V^*(s')]$$

This equation states that the value of a state under the optimal policy is the maximum expected return achievable by taking the best action in that state.

- **Action Value Function**:

$$Q^*(s, a) = \sum_{s'} P(s'|s, a)[R(s, a, s') + \gamma \max_{a'} Q^*(s', a')]$$

This equation states that the value of taking action $a$ in state $s$ is the immediate reward plus the discounted value of the best action in the next state $s'$.

The **Bellman equations** are fundamental to solving reinforcement learning problems because they provide a way to compute value functions and ultimately the optimal policy through **iterative methods** like **Dynamic Programming**, **Temporal-Difference (TD) Learning**, and **Q-learning**.

# Exploration vs. Exploitation

## 1. Exploration:

- **Exploration** refers to the agent's strategy of trying out new or unfamiliar actions to gather more information about the environment. The goal of exploration is to discover which actions lead to better rewards in states that have not yet been fully explored.

**Advantages:**

- **Learning Opportunities**: By exploring different actions, the agent can discover better strategies, even if the short-term reward is lower.
- **Avoiding Local Optima**: Exploration helps the agent avoid getting stuck in suboptimal solutions where a locally good action may prevent finding a globally better one.

**Challenges:**

- **Immediate Losses**: Exploration may lead to actions with low or even negative immediate rewards. This can be a costly process in terms of performance in the short term.
- **Risky Behavior**: Exploring unfamiliar actions can sometimes lead to poor outcomes, especially if the environment is stochastic or uncertain.

## 2. Exploitation:

- **Exploitation** refers to the agent's strategy of selecting actions based on the knowledge it has already gained to maximize immediate rewards. The agent leverages past experience to choose actions that are known to yield the highest reward in the current state.

**Advantages:**

- **Maximizing Immediate Gains**: Exploitation enables the agent to capitalize on what it already knows to achieve the best possible rewards in the short term.
- **Efficient Use of Knowledge**: By exploiting, the agent does not waste time on actions that it already knows are suboptimal.

**Challenges:**

- **Missed Opportunities**: Exploitation can prevent the agent from discovering better strategies or rewards if it stops exploring too early.
- **Stagnation**: If the agent focuses too much on exploitation, it may get stuck in a **local optimum**, where the strategy is good but not the best possible one.

## 3. Balancing Exploration and Exploitation:

The key challenge in RL is finding the right balance between exploration and exploitation. The agent needs to explore enough to gather sufficient knowledge about the environment but also exploit enough to maximize rewards. Various strategies help manage this balance:

**Approaches to Balance:**

1. **ε-greedy Strategy**:

- The **ε-greedy strategy** is a simple approach where the agent chooses a random action with probability $\varepsilon$ (exploration) and the best-known action with probability $1-\varepsilon$ (exploitation).
- Over time, the value of $\varepsilon$ can be decreased to reduce exploration as the agent becomes more confident in its knowledge of the environment.

2. **Decaying Exploration**:
   - In some algorithms, the agent starts with a high rate of exploration and gradually reduces exploration as it gathers more knowledge. This allows the agent to explore more early on and then exploit later.
3. **Boltzmann Exploration**:
   - In **Boltzmann exploration**, actions are selected based on a probability distribution that depends on their estimated values. Actions with higher rewards are chosen with higher probability, but less rewarding actions still have a chance of being selected. This introduces a more nuanced balance between exploration and exploitation.
4. **Upper Confidence Bound (UCB)**:
   - The **UCB algorithm** chooses actions based on the upper bound of the confidence interval for expected rewards. This means actions with less certainty (which could have high potential rewards) are explored more often, especially early in the learning process.

---

# Code Standards and Libraries Used in RL (Python/Keras/TensorFlow)

In reinforcement learning (RL), the choice of programming languages, libraries, and coding standards significantly impacts the development and efficiency of algorithms. Below is an overview of the primary coding standards and popular libraries used in RL, particularly focusing on **Python**, **Keras**, and **TensorFlow**.

## 1. Python:

- **Preferred Language**: Python is widely used in RL due to its simplicity, readability, and vast ecosystem of libraries.
- **Code Standards**:
  - **PEP 8**: The style guide for Python code promotes clean and readable code through consistent formatting and naming conventions.
  - **Docstrings**: Use docstrings for documenting functions and classes to improve code maintainability and clarity.
  - **Modularization**: Break code into reusable modules and functions to enhance reusability and testing.

## 2. Keras:

- **High-Level Library**: Keras is an open-source library built on top of TensorFlow that simplifies building and training neural networks.
- **Simplicity**: Keras provides a user-friendly API, allowing developers to quickly design, train, and evaluate models with minimal code.
- **Functional API**: Keras supports a functional API that enables the creation of complex models through layers, facilitating the implementation of custom architectures for RL.

**Key Features:**

- **Built-in Layers**: Predefined layers like Dense, Conv2D, and LSTM make it easy to build neural networks.
- **Callbacks**: Keras offers various callbacks for model monitoring and optimization, such as EarlyStopping and ModelCheckpoint.
- **Integration with TensorFlow**: Keras integrates seamlessly with TensorFlow, allowing for easy deployment and access to TensorFlow's functionalities.

## 3. TensorFlow:

- **Deep Learning Framework**: TensorFlow is a powerful open-source library for numerical computation and machine learning, widely used for building deep learning models, including those for RL.
- **Graph-Based Computation**: TensorFlow uses dataflow graphs, enabling efficient execution of operations on large datasets and deployment on various platforms (CPUs, GPUs, TPUs).

## Key Features:

- **TensorFlow Agents**: A specialized library for RL that provides implementations of various RL algorithms (like DQN, PPO, and A3C) and environments for testing.
- **Eager Execution**: This feature allows operations to execute immediately as they are called, making it easier to debug and experiment with models.
- **TensorFlow Serving**: Facilitates deploying trained models in production environments, allowing for scalable serving of ML models.

---

# Tabular Methods and Q-networks

## 1. Dynamic Programming (DP)

Dynamic Programming methods rely on a known model of the environment, including transition probabilities and reward structures. DP provides systematic ways to compute optimal policies and value functions.

## Key Concepts:

- **Policy Evaluation**: Assessing how good a policy is by calculating the expected return for each state.
- **Policy Improvement**: Adjusting the policy based on the value estimates to make it better.
- **Iterative Approach**: DP methods typically iterate between policy evaluation and policy improvement until convergence.

## Common DP Methods:

- **Policy Iteration**: Alternates between evaluating the current policy and improving it until an optimal policy is found.
- **Value Iteration**: Updates the value function for each state iteratively and derives the policy from the value function once convergence is achieved.

## Advantages:

- Guarantees convergence to the optimal policy when the model is known.
- Efficient for small state and action spaces.

**Limitations:**

- Requires complete knowledge of the environment model, which may not be available in real-world scenarios.
- Computationally expensive for large state spaces due to the exhaustive calculations involved.

## 2. Monte Carlo Methods

Monte Carlo methods estimate the value functions based on actual returns from sampled episodes of experience. These methods do not require a model of the environment and are useful for environments where the model is unknown.

### Key Concepts:

- **Episode-Based**: Monte Carlo methods rely on complete episodes to evaluate policies or state-action pairs.
- **Return Calculation**: The return for each state or action is calculated as the cumulative reward received from that point onward until the episode ends.

### Common Monte Carlo Methods:

- **First-Visit MC**: Estimates the value of a state by averaging the returns from the first time that state was visited in episodes.
- **Every-Visit MC**: Averages the returns from all occurrences of each state in the episodes.

### Advantages:

- Simple and intuitive, requiring minimal knowledge of the environment.
- Effective for episodic tasks where the termination is clear.

### Limitations:

- Requires a large number of episodes to converge to accurate value estimates.
- Cannot learn from incomplete episodes, limiting its applicability in continuous tasks.

## 3. Temporal-Difference Learning (TD Learning)

Temporal-Difference Learning combines ideas from both Monte Carlo methods and Dynamic Programming. It updates value estimates based on other learned estimates without requiring a model of the environment.

### *Key Concepts:*

- **Bootstrapping**: TD methods update the value function based on estimates from the current value function, which allows learning to occur on-line and in real time.
- **Value Updates**: TD learning updates value estimates based on the difference (temporal difference) between estimated values and observed rewards.

**Common TD Learning Methods:**

1. **TD(0):**

   - The simplest form of TD learning, where value estimates are updated based on the immediate reward and the estimated value of the next state.

   - Update rule:

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

   where $\alpha$ is the learning rate, $r_t$ is the reward, $\gamma$ is the discount factor, and $V(s)$ is the estimated value of state $s$.

2. **SARSA (State-Action-Reward-State-Action):**

   - An on-policy TD method that updates the action-value function based on the action taken and the subsequent state and action.

   - Update rule:

   $$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

3. **Q-Learning:**

   - An off-policy TD method that learns the optimal action-value function regardless of the policy being followed.

   - Update rule:

   $$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

   - Q-learning directly uses the maximum estimated future rewards to update the action-value function.

---

Deep Q-Networks (DQN) are a significant advancement in reinforcement learning, combining Q-learning with deep neural networks to handle high-dimensional state spaces. DQNs enable agents to learn optimal policies from raw sensory inputs, such as images, making them applicable in complex environments, such as playing video games.

---

# Deep Q-Networks (DQN)

**DQN** is an algorithm that utilizes a neural network to approximate the action-value function $Q(s,a)Q(s, a)Q(s,a)$, where sss is the state and aaa is the action. This approach addresses the limitations of traditional Q-learning, particularly in environments with large or continuous state spaces.

**Key Features:**

- **Function Approximation**: The neural network takes a state representation as input and outputs Q-values for all possible actions, allowing generalization across similar states.
- **Experience Replay**: DQNs store experiences in a replay buffer and sample mini-batches for training. This breaks the correlation between consecutive experiences, improving the stability and efficiency of learning.
- **Target Network**: DQNs maintain a separate target network, which is updated less frequently than the main Q-network. This helps stabilize training by reducing oscillations in the Q-value updates.

## DQN Algorithm:

1. Initialize the main Q-network and target network with random weights.

2. For each episode:

   - Observe the current state $s$.

   - Select an action $a$ using an exploration strategy (e.g., ε-greedy).

   - Execute action $a$ and observe the next state $s'$ and reward $r$.

   - Store the experience $(s, a, r, s')$ in the replay buffer.

   - Sample a mini-batch from the replay buffer.

   - Update the Q-network using the Bellman equation:
   $$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q_{\text{target}}(s', a') - Q(s, a)]$$

   - Periodically update the target network with the weights of the main Q-network.


## 2. Double Deep Q-Networks (DDQN)

**DDQN** addresses the overestimation bias often encountered in Q-learning. Traditional DQNs can sometimes overestimate Q-values due to using the same network for selecting and evaluating actions.

**Key Features:**

- **Action Selection and Evaluation**: In DDQN, the main network selects the action, while the target network evaluates the Q-value for that action. This decoupling helps reduce the overestimation bias.

## DDQN Algorithm:

1. The same initialization and experience replay mechanisms as in DQN are used.

2. During the Q-value update, use the main network to select the action:
$$a_{\text{max}} = \arg\max_a Q(s', a; \theta)$$
   and then use the target network to evaluate:

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q_{\text{target}}(s', a_{\text{max}}) - Q(s, a)]$$

## 3. Dueling DQNs

Dueling DQNs introduce a novel architecture that separates the representation of state values and advantages for different actions, allowing the agent to learn which states are valuable without requiring knowledge of the action values.

**Key Features:**

- **Value and Advantage Streams:** The network architecture has two separate streams: one for estimating the state value $V(s)$ and another for estimating the advantages $A(s, a)$ for each action. The Q-value is then computed as:

$$Q(s, a) = V(s) + (A(s, a) - \max_{a'} A(s, a'))$$

This formulation allows the network to learn the value of a state independent of the actions taken.

## 4. Prioritized Experience Replay

**Prioritized Experience Replay** enhances the experience replay mechanism by sampling experiences based on their significance, rather than uniformly. Experiences that are more informative for learning are prioritized, which can accelerate the learning process.

**Key Features:**

- **Priority Sampling**: Experiences are assigned a priority based on the magnitude of the TD error (the difference between predicted and actual returns). Higher TD errors indicate more informative experiences.
- **Importance Sampling**: To correct for the bias introduced by prioritized sampling, importance sampling weights are used during updates.

**Prioritized Experience Replay Algorithm:**

1. Store experiences in a replay buffer with associated priorities.
2. Sample mini-batches based on priority, with the most significant experiences being selected more frequently.
3. Update the Q-network using weighted updates based on the sampled experiences.

**Advantages of Prioritized Experience Replay:**

- Improves learning efficiency by focusing on experiences that contribute more to reducing the prediction error.
- Accelerates convergence and enhances overall performance of the agent.