

---

Certified Professional in AI Applications in Aviation

## Reinforcement Learning in Aviation Applications

---

**Action:** In Reinforcement Learning, an action refers to the decision made by an agent at each step of the learning process, based on the current state of the environment. Actions can be discrete (e.g., choosing among a set of predefined options) or continuous (e.g., selecting a value within a range).

**Agent:** An agent is the entity in Reinforcement Learning that interacts with the environment, making decisions (actions) and learning from the consequences (rewards or penalties). An agent can be a software program, a physical device, or even a human being.

**Convergence:** Convergence in Reinforcement Learning refers to the point at which the agent's learning process reaches a steady state, where the value function or policy no longer significantly changes. Convergence indicates that the agent has learned an optimal or near-optimal strategy for interacting with the environment.

**Deep Reinforcement Learning:** Deep Reinforcement Learning combines the principles of Reinforcement Learning with deep learning techniques, using artificial neural networks to approximate the value function or policy. Deep Reinforcement Learning enables agents to learn and adapt to complex, high-dimensional environments with many states and actions.

**Discount Factor:** The discount factor, denoted by  $\gamma$ , is a value between 0 and 1 that determines the significance of future rewards in Reinforcement Learning. A higher discount factor places more importance on future rewards, while a lower discount factor focuses on immediate rewards. The discount factor helps balance the trade-off between short-term and long-term rewards.

**Episode:** An episode is a sequence of interactions between an agent and the environment in Reinforcement Learning, starting from an initial state and continuing until a terminal state is reached. An episode may consist of multiple time steps, with the agent selecting actions and receiving rewards at each step.

**Markov Decision Process (MDP):** MDP is a mathematical framework used to model Reinforcement Learning problems. It describes the interaction between an agent and an environment with a set of states, actions, and transition probabilities between states. The Markov property in MDPs ensures that the future state depends only on the current state and action, and not on the history of past states and actions.

**Model-based Reinforcement Learning:** In Model-based Reinforcement Learning, the agent learns a model of the environment to predict the consequences of its actions. This model includes the transition probabilities between states and the rewards associated with each state-action pair. The agent then uses this model to plan its actions and optimize its strategy.

**Model-free Reinforcement Learning:** Model-free Reinforcement Learning does not rely on an explicit model of the environment. Instead, the agent learns the value function or policy directly from interactions with the

environment. Model-free methods are often sample-efficient and can handle complex environments without requiring an accurate model.

**Multi-agent Reinforcement Learning (MARL):** MARL is a subfield of Reinforcement Learning where multiple agents interact and learn concurrently in a shared environment. MARL can be used to model cooperative, competitive, or mixed settings, where agents may need to coordinate, compete, or negotiate with each other to achieve their goals.

**Policy:** A policy is a mapping from states to actions that defines the agent's behavior in a Reinforcement Learning problem. It specifies which action the agent should take given a particular state, and can be deterministic (selecting a single action) or stochastic (assigning probabilities to actions).

**Q-function (Q-value):** The Q-function, also known as the action-value function, estimates the expected cumulative reward of taking a specific action in a given state and following a particular policy thereafter. The Q-function helps the agent determine the best action to take at each time step.

**Q-learning:** Q-learning is a popular Reinforcement Learning algorithm that learns the Q-function for a given Markov Decision Process. Q-learning updates the Q-values iteratively, using temporal difference methods, until convergence to the optimal Q-function.

**Reinforcement Learning:** Reinforcement Learning is a subfield of machine learning that focuses on training agents to make decisions in complex, dynamic environments through trial and error. Reinforcement Learning agents learn by interacting with the environment, receiving rewards or penalties for their actions, and updating their strategies based on these experiences.

**Reward:** A reward is a scalar value that the agent receives upon transitioning from one state to another in the environment. Rewards serve as feedback to the agent, guiding its learning process towards optimal behavior. Rewards can be positive (indicating success or progress) or negative (indicating failure or setbacks).

**State:** A state is a representation of the environment at a particular time step in Reinforcement Learning. It encapsulates all the relevant information that the agent needs to make decisions and take actions. States can be discrete (e.g., a set of predefined categories) or continuous (e.g., a vector of real-valued variables).

**State-action Value Function:** The state-action value function, also known as the Q-function, estimates the expected cumulative reward of taking a specific action in a given state and following a particular policy thereafter. It helps the agent determine the best action to take at each time step.

**Temporal Difference Learning (TD Learning):** TD Learning is a family of Reinforcement Learning algorithms that update the value function based on the observed difference between predicted and actual rewards. TD Learning combines the advantages of dynamic programming and Monte Carlo methods, enabling agents to learn from delayed rewards and avoid the need for a full episode's trajectory.

**Value Function:** The value function, also known as the state-value function, estimates the expected cumulative reward of being in a particular state and following a specific policy. It helps the agent evaluate

the desirability of different states and guide its learning process towards optimal behavior.

**Value Iteration:** Value Iteration is a Reinforcement Learning algorithm that computes the optimal value function iteratively, starting from an initial guess and updating it based on the Bellman optimality equation until convergence. Once the optimal value function is obtained, the agent can derive the optimal policy by selecting the action with the highest Q-value for each state.