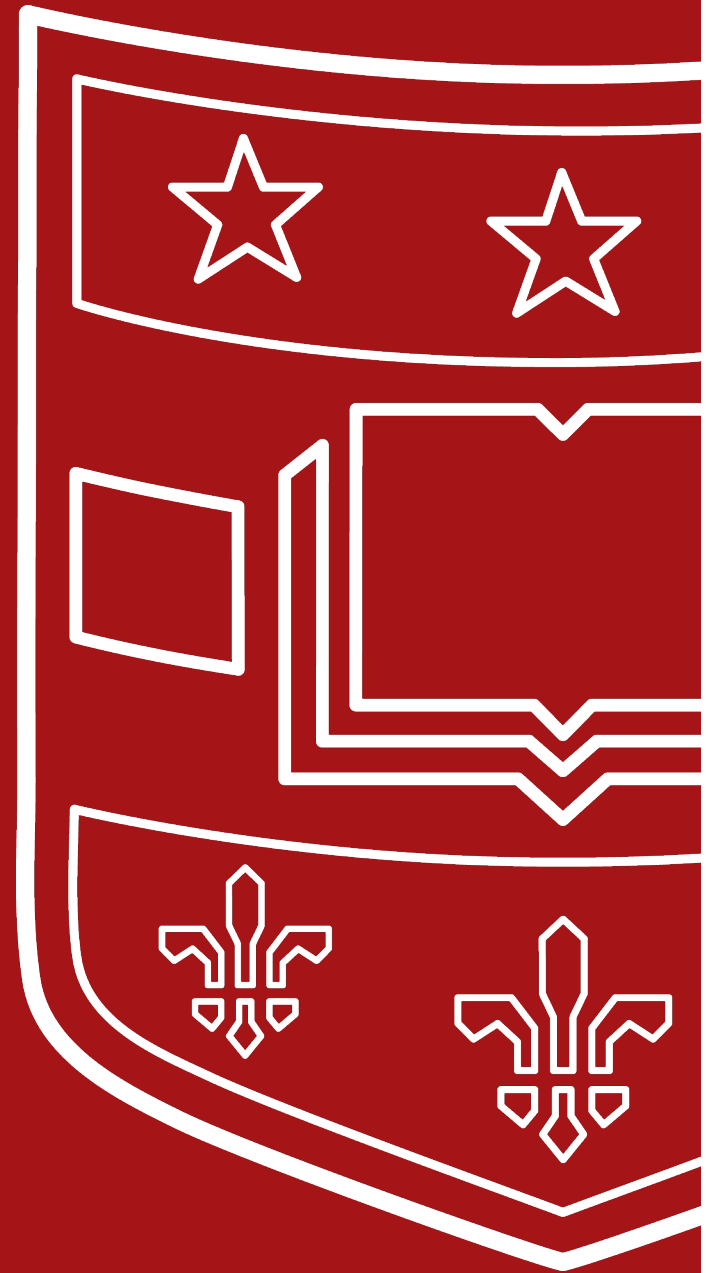


# Introduction to Reinforcement Learning

Ruopeng An





# What is Reinforcement Learning?

- Definition
  - Reinforcement Learning is a machine learning method where an agent learns by interacting with an environment to achieve goals.
- Key Objective
  - The primary goal in Reinforcement Learning is to develop a policy that maximizes cumulative rewards over time.
- Unique Aspect
  - Reinforcement Learning stands out for its focus on learning through interaction without explicit task programming.



# Key Concepts of RL

- Agent
  - The agent in RL is the learner or decision maker that interacts with the environment to achieve goals.
  - It is responsible for selecting actions based on the current state to maximize rewards.
- Environment
  - The environment in RL is where the agent operates, including all elements it interacts with and learns from.
  - It encompasses the context in which the agent makes decisions and receives feedback through actions.
- Actions
  - Actions in RL represent the set of possible moves the agent can make in the environment.
  - They are the decisions taken by the agent to influence the state and receive rewards or penalties.
- States
  - States in RL describe the current situation returned by the environment to the agent.
  - They provide information about the context in which the agent is operating, influencing its decision-making process.



# Key Concepts of RL

- Rewards
  - Rewards in RL are feedback from the environment that indicates the value of actions taken by the agent.
  - They serve as positive reinforcement for actions that lead to desirable outcomes.
- Penalties
  - Penalties in RL are negative rewards given to discourage certain actions by the agent.
  - They help the agent learn to avoid actions that lead to unfavorable outcomes.
- Policy
  - The policy in RL is the strategy the agent uses to determine the next action based on the current state.
  - It guides the agent's decision-making process by mapping states to actions to achieve long-term goals.
- Value Function
  - The value function in RL estimates the expected return of states or state-action pairs.
  - It helps the agent evaluate the potential outcomes of different actions in a given state.
- Discount Factor
  - The discount factor in RL measures how much future rewards are worth compared to immediate rewards.
  - It influences the agent's decision-making by balancing short-term gains with long-term benefits.



# Advantages of RL

- **Adaptability**
  - Reinforcement Learning showcases remarkable adaptability by learning and adjusting strategies based on environmental feedback.
  - The agent's ability to adapt to dynamic environments sets RL apart in problem-solving scenarios.
- **Decision Making**
  - RL excels in decision-making by strategically selecting actions to achieve long-term objectives.
  - The agent's sequential decision-making process in RL leads to optimal strategies for goal attainment.
- **Complex Problem Solving**
  - RL tackles complex, multi-step challenges effectively by breaking them down into manageable actions.
  - The methodical approach of RL enables it to navigate intricate problems with hidden solutions.
- **Examples**
  - RL applications span diverse fields like robotics, game playing, and autonomous vehicles, showcasing its versatility.
  - From chess-playing algorithms to self-driving cars, RL demonstrates its effectiveness in real-world applications.



# Limitations of RL

- **Sample Efficiency**
  - Reinforcement Learning requires extensive data for effective learning, demanding numerous experiences to refine decision-making processes.
  - The need for substantial data sets can pose challenges in real-world applications, hindering quick adaptation.
- **Complexity**
  - Designing the reward system and state representation in RL can be intricate, requiring careful consideration of various factors.
  - The complexity arises from creating a system that accurately evaluates actions and states, impacting the learning process.
- **Computationally Intensive**
  - RL often demands significant computational resources, especially in high-dimensional spaces, for efficient learning and decision-making.
  - The computational intensity stems from processing large amounts of data and complex algorithms, necessitating powerful computing capabilities.



# RL vs. Other Machine Learning Approaches

- Supervised Learning
  - In Supervised Learning, models learn from labeled data to predict outcomes accurately.
  - It relies on labeled examples to train algorithms for classification or regression tasks.
- Unsupervised Learning
  - Unsupervised Learning focuses on finding patterns or structures in data without labeled examples.
  - It aims to explore and uncover hidden relationships or structures within the data.
- Evolutionary Algorithms
  - Evolutionary Algorithms simulate evolution to optimize solutions over generations.
  - They use mechanisms like mutation and selection to evolve solutions towards an optimal outcome.



# The Agent in RL

- Definition
  - In Reinforcement Learning, the agent is the learner that interacts with the environment to achieve a goal.
  - The agent in RL is the entity that makes decisions based on interactions with the environment.
- Example
  - An example of an agent in RL is a chess-playing program that learns to make moves to win the game.
  - Consider a robot navigating a maze as an example of an agent in Reinforcement Learning.





# Episodic vs. Continuous Agents

- Episodic Agent
  - Operates in separate episodes with terminal states for each interaction.
  - Well-suited for tasks with clear breaks and distinct episodes.
- Continuous Agent
  - Operates in a continuous flow without distinct episode boundaries.
  - Ideal for tasks where actions and states seamlessly transition.
- Example
  - Puzzle-solving agent with each puzzle as an episode.
  - Temperature control system adjusting continuously in a building.



# Deterministic vs. Stochastic Agents

- Deterministic Agent
  - A deterministic agent always produces the same outcome in a given state.
  - Examples include calculators where specific inputs yield consistent results.
- Stochastic Agent
  - Stochastic agents introduce uncertainty as outcomes can vary in the same state.
  - An example is a stock trading agent affected by market volatility.



# Understanding the Environment

- Example
  - In the context of an autonomous vehicle, the environment includes roads, traffic signals, vehicles, pedestrians, and obstacles, influencing the vehicle's decision-making.



# State Representation and Action Selection

- State Representation
  - Describes the current situation perceived by the agent from the environment.
  - Encodes essential information about the environment for decision-making purposes.
- Action Selection
  - Determines the optimal action for the agent to take based on the current state.
  - Involves choosing from available actions to maximize rewards or achieve goals.
- Example
  - In a video game, state representation could include player positions and object locations.
  - Action selection in the game might involve moving, jumping, or interacting with objects based on the current state.



# Action Selection Mechanisms

- Random Exploration
  - Random exploration involves selecting actions without considering the current state or policy.
  - This method allows the agent to discover new possibilities and learn about the environment through trial and error.
- Epsilon-Greedy
  - Epsilon-Greedy is a strategy that balances exploration and exploitation in action selection.
  - It involves choosing a random action with a certain probability (epsilon) while favoring known best actions with the remaining probability.
- Policy-Based Methods
  - Policy-based methods select actions based on a learned policy that maps states to actions.
  - These methods leverage past experiences to make decisions, improving performance over time.



# Introducing Markov Decision Processes (MDPs)

- Definition: A Markov Decision Process is a mathematical framework used to model decision-making in situations where outcomes are partly random and partly under the control of a decision-maker.
- Components:
  - States (S): The set of all possible states in the environment.
  - Actions (A): The set of all actions the agent can take.
  - Transition Probability (P):  $P(s'|s, a)$ , the probability of transitioning from state  $s$  to state  $s'$  after taking action  $a$ .
  - Rewards (R):  $R(s, a, s')$ , the reward received after transitioning from  $s$  to  $s'$ , due to action  $a$ .
  - Discount Factor ( $\gamma$ ): Represents the importance of future rewards.



# The Significance of MDPs in RL

- Foundation for Many RL Problems: MDPs provide the theoretical underpinnings for most reinforcement learning problems, offering a formalization for environments in which an agent operates.
- Solving MDPs: Solving an MDP involves finding a policy ( $\pi$ ) that maximizes some measure of long-term reward. Algorithms like Value Iteration and Policy Iteration are used to find optimal policies.
- Real-World Applications: MDPs are used to model and solve a wide range of real-world problems, from robot navigation and automated control to economic decision-making and game strategy development.



# Introduction to the Bellman Equation

- Objective: Understand the Bellman equation's role in reinforcement learning.
- Overview: The Bellman equation provides a recursive solution for finding the optimal policy in both deterministic and stochastic environments.





# Bellman Equation - Deterministic Environments

- Definition: In deterministic environments, the outcome of each action is predictable.
- Equation:  $V(s) = \max_a [R(s, a) + \gamma V(s')]$ 
  - $V(s)$ : Value of state  $s$ .
  - $\max_a$ : Maximizing over all possible actions  $a$ .
  - $R(s, a)$ : Reward received after taking action  $a$  in state  $s$ .
  - $\gamma$ : Discount factor.
  - $V(s')$ : Value of the resulting state  $s'$ .
- Explanation: The equation calculates the value of each state by considering the immediate reward plus the discounted value of the next state.



# Example - Deterministic Environment

- Context: A simple grid-world game where an agent moves deterministically.
- Application: Use the Bellman equation to calculate the value of each grid position, aiding the agent in finding the optimal path to the goal.



# Bellman Equation - Stochastic Environments

- Definition: In stochastic environments, the outcome of actions is not guaranteed, introducing probability into the decision process.
- Equation:  $V(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')]$
- $P(s'|s, a)$ : Probability of transitioning to state  $s'$  from state  $s$  after taking action  $a$ .
- $R(s, a, s')$ : Reward received for transitioning from  $s$  to  $s'$  due to action  $a$ .
- Explanation: The value of each state considers all possible next states, weighted by their probability of occurrence and potential rewards.



## Example - Stochastic Environment

- Context: A board game where dice rolls determine the agent's movement.
- Application: Use the Bellman equation to calculate the expected value of each board position, considering all possible outcomes of the dice roll.



# Solving Problems with the Bellman Equation

- Dynamic Programming: Techniques like Value Iteration and Policy Iteration can solve the Bellman equation, finding the optimal policy for both deterministic and stochastic environments.
- Significance: Understanding the Bellman equation is crucial for developing effective reinforcement learning algorithms that can navigate the complexities of real-world decision-making.



# Introduction to Q-Learning

- Objective: Understand Q-learning, a model-free reinforcement learning algorithm.
- Overview: Q-learning enables an agent to learn the value of an action in a particular state, guiding it toward the goal without a model of the environment.



# What is Q-Learning?

- Definition: Q-learning is a form of model-free reinforcement learning that does not require knowledge of the environment's dynamics. Instead, it learns the value of action-state pairs by exploration.
- Goal: The primary goal is to learn the policy that maximizes the total reward.



# The Q-Learning Equation

- Equation:  $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$ 
  - $Q(s, a)$ : The current estimate of the action value.
  - $\alpha$ : Learning rate.
  - $r$ : Reward received after taking action  $a$  in state  $s$ .
  - $\gamma$ : Discount factor, indicating the importance of future rewards.
  - $\max_{a'} Q(s', a')$ : The maximum predicted reward, obtained from the next state  $s'$ .
- Explanation: The Q-learning equation updates the Q-value of the state-action pair, incorporating the immediate reward plus the discounted maximum future reward.





# Key Components of Q-Learning

- -Learning Rate ( $\alpha$ ): Determines how much new information overrides old information.
- Discount Factor ( $\gamma$ ): Balances the importance of immediate and future rewards.
- Policy: Derived from Q-values, often using an  $\epsilon$  –greedy strategy for exploration and exploitation.



# Estimating $\max_{a'} Q(s', a')$

- Objective: To understand how the future reward is estimated in Q-learning.
- Estimation of Future Reward:
  - The term  $\max_{a'} Q(s', a')$  represents the estimated maximum future reward obtainable from the next state  $s'$ , considering all possible actions  $a'$ .
  - This estimation is crucial for the algorithm because it provides a lookahead feature, enabling the agent to evaluate the potential future benefits of its current actions.



# How to Estimate $\max_{a'} Q(s', a')$

- Exploration and Learning:
  - Initially,  $Q(s', a')$  values are unknown and typically initialized to a default value (e.g., zeros).
  - As the agent explores the environment and receives rewards for its actions, it updates the  $Q$  values based on actual outcomes.
- Update Process:
  - The agent updates its knowledge about the value of taking action  $a$  in state  $s$  towards optimizing future rewards.
  - Through repeated interactions with the environment, the agent refines its estimates of  $Q(s', a')$ , leading to more accurate predictions of future rewards.

# Purpose of Subtracting $Q(s, a)$ in the Update Equation



- Objective: Understand the role of subtracting  $Q(s, a)$  in the Q-learning update rule.
- Temporal Difference Error:
  - The term  $[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$  is known as the temporal difference (TD) error.
  - This error represents the difference between the agent's current estimate of the future reward ( $Q(s, a)$ ) and the newly observed estimate ( $r + \gamma \max_{a'} Q(s', a')$ ).



# Significance of the Temporal Difference Error

- Learning from Mistakes:
  - The subtraction of  $Q(s, a)$  allows the algorithm to adjust the value of  $Q(s, a)$  closer to the newly observed estimate, correcting overestimations or underestimations.
- Convergence to Optimal Policy:
  - By continuously updating the  $Q$  values based on the TD error, the algorithm incrementally learns the optimal policy that maximizes the cumulative reward.
- Adaptation and Improvement:
  - This mechanism enables the agent to adapt its strategy based on its experiences, refining its approach as it learns more about the environment.



# Conclusion on Q-Learning Update Dynamics

- Recap: The estimation of  $\max_{a'} Q(s', a')$  and the subtraction of  $Q(s, a)$  are essential components of the Q-learning algorithm, allowing it to predict future rewards and adjust its strategy based on learning from the temporal difference error.
- Takeaway: These mechanisms are fundamental to the algorithm's ability to learn optimal actions for maximizing rewards in various environments.



# Example - Navigating a Maze

- Context: An agent learns to navigate through a maze to reach a goal.
- Application: Q-learning can be applied without a model of the maze. The agent explores the maze, updating Q-values based on the rewards received for reaching the goal.



# Implementing Q-Learning

- Initialization: Start with arbitrary Q-values.
- Loop: For each state-action pair, update Q-values using the Q-learning equation.
- Convergence: With enough exploration, Q-values converge to optimal values, allowing the agent to follow the optimal policy.





# Implementing Q-Learning

- Initialization: Start with arbitrary Q-values.
- Loop: For each state-action pair, update Q-values using the Q-learning equation.
- Convergence: With enough exploration, Q-values converge to optimal values, allowing the agent to follow the optimal policy.



# Advantages and Challenges of Q-Learning

- Advantages:
  - Model-free: Does not require a model of the environment.
  - Flexibility: Applicable to a wide range of problems.
- Challenges:
  - Requires significant exploration.
  - May converge slowly in complex environments.