

# Chapter XVI - Making Complex Decisions

## 16.1 Sequential Decision Problems

**Sequential decision problems** involve making a series of decisions over time, where the outcome of one decision influences future decisions. These problems are common in many real-world applications, such as robotics, finance, and healthcare, where the agent must choose actions not just based on the current state but also considering the effects on future states.

### Key Concepts:

- **States:** Represent the current condition of the system or environment.
- **Actions:** The choices or decisions that an agent can make.
- **Policies:** A strategy that defines which action to take in each state.
- **Rewards:** The feedback the agent receives from the environment after taking an action in a given state.
- **Transition Model:** Describes the probability of transitioning from one state to another given an action.

### Example:

In a robotic vacuum cleaner, the decision to move to a certain area of the room depends on the robot's current state (location) and the subsequent cleaning requirements (future states), which could affect the next actions (e.g., charging when the battery is low).

---

## 16.2 Value Iteration

**Value Iteration** is a method for solving sequential decision problems by calculating the value of each state. The goal is to find the optimal policy that maximizes the expected cumulative reward (also known as the **value function**) for an agent.

### Key Concepts:

- **Value Function ( $V(s)$ ):** Represents the maximum expected cumulative reward achievable from state  $s$ .
- **Bellman Equation:** A recursive relationship used to calculate the value function. It relates the value of a state to the values of its possible successor states, considering the reward

and the transition model.

Value Iteration works by iteratively updating the value of each state until convergence, at which point the optimal policy can be derived by selecting the action that leads to the highest value in the next state.

### Steps in Value Iteration:

1. Initialize the value function for all states.
2. For each state, update its value using the Bellman equation.
3. Repeat until the value function converges.

### Example:

In a gridworld environment, an agent may have several possible movements (up, down, left, right) from each state. The value of each state is computed by considering the rewards of all possible future states.

---

## 16.3 Policy Iteration

**Policy Iteration** is another approach to solving sequential decision problems, where the goal is to improve an initial policy iteratively. It works by alternating between **policy evaluation** (assessing how good the current policy is) and **policy improvement** (updating the policy based on the evaluation).

### Key Concepts:

- **Policy:** A mapping from states to actions.
- **Policy Evaluation:** Calculates the value function under the current policy.
- **Policy Improvement:** Updates the policy by choosing actions that maximize the expected value of the future states.

### Steps in Policy Iteration:

1. **Initialize a policy** for each state.
2. **Evaluate the policy** by computing the value function of all states based on the current policy.
3. **Improve the policy** by selecting the action that maximizes the value function in each state.
4. Repeat the evaluation and improvement steps until the policy stabilizes.

## Example:

In a maze-solving problem, the agent starts with an arbitrary policy (random actions) and iteratively refines it by evaluating the best actions to take at each position in the maze based on the expected rewards.

---

## 16.4 Partially Observable MDPs (POMDPs)

A **Partially Observable Markov Decision Process (POMDP)** extends the traditional Markov Decision Process (MDP) by allowing the agent to have **partial observability** of the environment. This means the agent does not have full information about the current state, but only receives **observations** that provide partial information about the true state.

### Key Concepts:

- **Belief State:** A probability distribution over all possible states, representing the agent's belief about the actual state based on observations.
- **Observation Function:** Describes the probability of receiving an observation given a state and action.
- **Action-Observation Cycle:** The agent takes an action, receives an observation, and updates its belief state accordingly.

Solving POMDPs involves computing an optimal policy based on belief states, which can be more computationally challenging than in fully observable MDPs.

## Example:

Consider a self-driving car navigating through foggy conditions. The car cannot directly observe its exact location but receives sensor data that gives partial information about its surroundings (e.g., proximity to other vehicles).

---

## 16.5 Multi-Agent Decisions: Game Theory

**Game Theory** is the study of mathematical models of strategic interactions between rational agents. In multi-agent decision problems, each agent's decisions affect the outcomes for other agents, and they may have conflicting or cooperative goals.

### Key Concepts:

- **Players:** The agents involved in the game.
- **Strategies:** The plans of action each player may take.
- **Payoff Functions:** Represent the utility or benefit a player gains from a particular outcome of the game.
- **Nash Equilibrium:** A set of strategies where no player can improve their payoff by unilaterally changing their strategy.

Game theory is applied to various domains, including economics, politics, and artificial intelligence, where agents (e.g., competitors, collaborators) interact and make decisions that depend on each other's choices.

### Example:

In a pricing game, multiple companies may adjust their prices to maximize profits, knowing that each company's price influences the others. The equilibrium state occurs when no company wants to change their price, given the prices of the others.

---

## 16.6 Mechanism Design

**Mechanism Design** is a branch of game theory concerned with designing rules and systems that lead to desired outcomes, even when participants may act in their own self-interest. The goal is to **incentivize** agents to behave in a way that achieves a socially optimal outcome.

### Key Concepts:

- **Incentive Compatibility:** The property that ensures agents are motivated to reveal their true preferences or information.
- **Social Welfare:** The overall well-being or utility of the group of agents.
- **Auction Design:** A common application of mechanism design, where rules are created to allocate goods or resources efficiently.

Mechanism design is widely used in economics, auctions, and the allocation of resources in multi-agent systems.

### Example:

A central auction system for allocating goods (like electricity or bandwidth) can be designed to encourage truthful bidding. In such a system, participants are incentivized to bid their true value of the goods, leading to efficient resource allocation.

---

# Exercises

1. **Value Iteration Exercise:** Given a gridworld environment with rewards at specific locations (e.g., -1 for obstacles, +10 for goals), implement value iteration to find the optimal policy that maximizes the expected reward for the agent starting at a given location.
2. **Policy Iteration Exercise:** Implement policy iteration for a simple decision problem where an agent must decide whether to invest in a project with uncertain outcomes (success or failure). The agent should maximize its expected utility.
3. **POMDP Exercise:** Simulate a robot in a partially observable environment (e.g., a maze with some blocked paths). Use belief states to update the robot's knowledge of its location and actions to find the best path to the goal.
4. **Game Theory Exercise:** Model a two-player game where each player has two strategies: "Cooperate" or "Defect." Create a payoff matrix and find the Nash equilibrium.
5. **Mechanism Design Exercise:** Design an auction mechanism for allocating limited resources, such as bandwidth or cloud storage, that incentivizes participants to bid truthfully about their preferences.