

The background is a dark blue gradient. It features several decorative elements: concentric circles in the top-left and bottom-right corners, and various squares of different sizes scattered throughout. Some squares are outlined in a light teal color, while others are white. The text is centered and reads:

Extended Essay Sample

# MATHEMATICS

Mathematical concepts of MDP and Graphs are  
used in Machine Learning for Robots



## Alright, DP legends, ready to own that next assignment?

Pop your specs into our **order form**—roll with a similar draft or a fresh custom—and punch in the coupon code **STARTDP20** for **20% OFF** your first banger.

We're talkin' total rubric lockdown, IB-criteria flex, and spot-on instructions.

Sit back, chill, and watch those marks glow up 👍

**ORDER NOW** →



**Title: Mathematical concepts of MDP and Graphs are used in Machine Learning for Robots**

**Word Count: 6517**

**Your Name:**

**Candidate Number:**

# **Mathematical concepts of MDP and Graphs are used in Machine Learning for Robots**

## **Introduction**

The blending of mathematical ideas with machine learning has opened the way for tremendous developments, notably in the realm of robotics, in the ever-evolving environment of technology and artificial intelligence. The goal of this extended essay is to go deeply into the intricate mathematics that supports the confluence of graph theory and Markov Decision Processes (MDP) in the context of machine learning for robotics.

Robotics is being used more and more often in a variety of fields and applications, from autonomous cars and industrial automation to healthcare and space exploration. These robots have developed into sentient beings who can make judgments and learn from their environment. They are no longer only mechanical machines. The integration of machine learning and mathematics allows for this transition.

The mathematical framework known as Markov Decision Processes, or MDPs, is used to describe decision-making issues where results are partially random and partially within the control of the decision-maker. In the context of robotics, where judgments must be made autonomously while navigating dynamic and unpredictable settings, this formalism is essential. We learn how robots may make decisions that maximize rewards and minimize dangers by examining MDPs.

On the other hand, graph theory offers a flexible set of mathematical tools for simulating and examining the connections and interconnections in complicated systems. Graphs are used to depict geographical contexts, sensor networks, and communication systems in the context of robotics. Robots are better able to travel, communicate, and coordinate tasks when they are aware of how to use graph theory.

## **Aim**

I hope to reveal complex relationships between MDPs, Graph Theory, and machine learning in the context of robotics by completing this lengthy article. My objective is to demonstrate not just their theoretical foundations but also their practical relevance and influence. In addition, I want to provide readers with a nuanced knowledge of the sophisticated mathematical ideas that are incorporated into this multidisciplinary topic.

I will make my way through the mathematical terrains of MDPs and Graph Theory, investigate how they work in conjunction with machine learning methods, and clarify how they improve the capabilities of robotic systems. I wish to contribute to a better understanding of how mathematics is influencing the future of robots and, by extension, our world through careful study, specific examples, and useful ideas.

## **Background**

### **Reinforcement Learning:**

Machine learning's reinforcement learning (RL) area focuses on creating algorithms that can make judgments sequentially to maximize cumulative rewards in unpredictable environments. It is a mathematical framework that enables agents, like robots, to pick up knowledge by interacting with their environment, much like people do through trial and error.

### **Graph Theory**

The study of graphs, which are mathematical structures made up of nodes (vertices) and edges (connections between them), is one of the topics covered by the field of mathematics known as graph theory. From social networks and transportation systems to computer networks and physical contexts, graphs are used to describe interactions, connectedness, and networks in a variety of fields.

Nodes, edges, routes, and cycles are the basic building blocks of graph theory. Graphs can be either undirected (where the edges have no direction) or directed (where the edges have a direction). They can also be unweighted or weighted (where edges have associated values).

### **Relationship and goals**

When used in the context of robotics, RL and Graph Theory are related. Robots frequently work in dynamic and spatially complicated situations, which need critical navigation and decision-making skills. In such circumstances, the robot's decision-making process may be modeled using MDPs, a fundamental idea in RL. The spatial organization of the environment may be represented by graphs from the field of graph theory, allowing for effective path planning and resource allocation for the robot.

The goal of graph theory is to examine and comprehend the characteristics and structure of graphs. This comprises, among other things, locating the shortest pathways, figuring out connectedness, locating clusters, and enhancing network flow. Graph theory is frequently used in robotics to aid effective resource allocation, communication, and navigation.

To support multi-agent decision-making in mixed autonomy traffic, autonomous mobility-on-demand (AMoD) systems, and autonomous exploration, researchers have developed a variety of graph reinforcement learning (GRL) methodologies. These methods are successful in improving the efficiency of multi-agent decision-making, making it possible to recover transferable, generalizable, and scalable behavior regulations, and predicting a robot's best sensing course of action in belief space.

Q-learning and Deep Q-Networks (DQN) are two reliable reinforcement learning algorithms that have found wide use in robotics.

### **Q-Learning:**

Q-learning is a fundamental reinforcement learning (RL) technique that works with Markov Decision Processes (MDPs). It is founded on the idea of estimating a Q-function, represented as  $Q(s, a)$ , which stands for the expected cumulative benefit of adopting action "a" in the state "s" and then adhering to an ideal course of action thereafter. The Bellman equation, which mathematically describes the value of  $Q(s, a)$  as a recursive connection involving the maximum Q-values of subsequent states, is the fundamental update rule of Q-learning. The theory of stochastic processes and dynamic programming form the mathematical basis of Q-learning.

### **Deep Deterministic Policy Gradient (DDPG):**

For continuous action spaces, the value-based reinforcement learning algorithm DDPG was developed. It combines the advantages of Q-learning with policy gradients to deal with the complicated control issues that robots frequently face.

For continuous action spaces, the value-based reinforcement learning algorithm DDPG was developed. It combines the advantages of Q-learning with policy gradients to deal with the complicated control issues that robots frequently face.

### **Challenges:**

Exploration vs. exploitation:

It might be difficult mathematically to strike a balance between exploitation (selecting actions with the largest anticipated rewards) and exploration (trying different actions for learning). DDPG frequently injects noise into gameplay to promote exploration.

Function Approximation:

When deep neural networks are used to approximate Q-values and policy functions, problems with approximation stability and convergence arise. To deal with these problems, strategies like experience replay and target networks are used.

## **Methods and techniques**

The following Methods and techniques were employed to investigate the function of graph theory in reinforcement learning for autonomous robots:

### **1. Graph-Based State Representation**

This approach uses a graph to describe the state space of the environment. While edges signify changes in state depending on robot activities, nodes indicate various states or setups. Researchers can effectively model the state space and adapt it for reinforcement learning algorithms by using Graph Theory principles.

### **2. Graph-Based Reward Shaping**

In this method, reward functions are defined using a graph. State nodes and state transitions are represented by graph nodes and edges, respectively. To direct a robot's learning process, researchers give incentives to states or edges. The reinforcement learning agent receives more organized input thanks to graph-based reward shaping.

### **3. Graph-Based Path Planning**

A navigation graph depicting the robot's surroundings is built using graph theory. In this network, the best pathways are found using algorithms like Dijkstra's. Agents using reinforcement learning may then develop rules based on these pathways, making it possible for them to navigate complicated environments with effectiveness and efficiency.

## **Mathematical model of Markov Decision Processes (MDPs)**

Markov Decision Processes (MDPs) in Reinforcement Learning for Autonomous Robots: A Mathematical Model

### **MDPs Distinct Concepts:**

1. Transition probabilities ( $p$ ): The transition probabilities, abbreviated as  $P$  in an MDP, specify the probability distribution across potential future states given the present state and action.

2. Reward Function (R): The reward function, designated as R, connects each state-action pair with a numerical reward. It may be expressed as  $R(s, a, s')$ , which denotes the reward obtained after changing from state  $s$  to state  $s'$  by acting  $a$ .
3. Policy( $\pi$ ): A policy, represented by the symbol, represents a mapping from states to actions and reveals the agent's method for deciding which actions to do in each state. It outlines how the agent will act in its surroundings.
4. Value Function (V): The value function, symbolized by the letter V, depicts the anticipated cumulative benefit an agent can obtain from a specific condition while adhering to a specific policy.

## **Mathematical model of MDP**

### **Problem Statement:**

Consider a situation where an autonomous robot must find its way from point A to point E in a graphed environment. The goal of the robot is to choose the best path while considering the rewards and transition probabilities of various actions.

### **Mathematical Exploration:**

#### **Numerical Example 1:**

##### **Step 1: Problem Setup**

We are trying to determine the best route from A to E in a graph that represents the environment and has the nodes A, B, C, D, and E. The distances between the nodes are known.

##### **Step 2: Graph State Representation**

Represent the environment as a graph, where nodes represent locations, and edges represent possible transitions between locations. We'll assign weights to the edges, denoting distances or costs associated with transitions.

$$\text{Distance}(A, B) = 2$$

$$\text{Distance}(A, C) = 4$$



$$\text{Distance}(B, C) = 1$$

$$\text{Distance}(B, D) = 7$$

$$\text{Distance}(C, D) = 3$$

$$\text{Distance}(C, E) = 5$$

$$\text{Distance}(D, E) = 2$$

### Step 3: Markov Decision Processes (MDPs)

Define the MDP components:

States (S): The states represent grid cells in the environment.  $S = \{0, 1, 2, 3, 4, 5\}$  (inventory levels).

- Actions (A): The robot can take actions "up," "down," "left," "right," or "stay."  $A = \{\text{"up"}, \text{"down"}, \text{"left"}, \text{"right"}, \text{"stay"}\}$  or  $A = \{0, 1, 2, 3\}$ .
- Transitions (T): Transition probabilities represent the likelihood of moving from one state to another based on the chosen action.

### Step 4: Transition Probabilities Calculation

For this step, we'll calculate the transition probabilities for state-action pairs based on the distances between nodes. We'll use the logistic function formula you provided:

$$P(A \rightarrow B) = 1 / (1 + e^{(\text{Distance}(A, B) - \text{Mean}(\text{distance})))}$$

Let's calculate the transition probability for one example:

For state A (current location) and action "move to adjacent node B" (transition to B):

$$\text{Distance}(A, B) = 2 \text{ (given)}$$

Now, calculate the transition probability using the logistic function formula:

$$P(A \rightarrow B) = 1 / (1 + e^{(2 - \text{Mean}(\text{distance}))})$$

Assuming  $\text{Mean}(\text{distance})$  is some constant (let's say 3 for this example):

$$P(A \rightarrow B) = 1 / (1 + e^{(2 - 3)})$$

$$P(A \rightarrow B) = 1 / (1 + e^{(-1)})$$

$$P(A \rightarrow B) \approx 0.731$$

Repeat this calculation for all state-action pairs to determine their transition probabilities.

### Step 5: Value Iteration

In this step, we'll iteratively update the value function  $V(s)$  using the Bellman equation:

$$V(s) = \max[\sum(P(s' | s, a) \times (R(s, a) + \gamma \times V(s')))]$$

Assume a discount factor ( $\gamma$ ) of 0.9.

Let's calculate  $V(s)$  for one example:

#### For state A:

Initialize  $V(A) = 0$  (initial assumption).

For action "move to adjacent node B" (transition to B):

Calculate the expected value (weighted sum of the next states' values):

$$V(A) = \max[P(A \rightarrow B) \times (R(A, B) + 0.9 \times V(B))] \text{ (using the transition probability and reward for this action)}$$

Assuming we've already calculated  $P(A \rightarrow B)$  as approximately 0.731 and the reward  $R(A, B)$  is -2 (negative of distance), and let's assume  $V(B)$  is 0 (initial assumption):

$$V(A) = \max[0.731 \times (-2 + 0.9 \times 0)]$$

$$V(A) = \max[0.731 \times (-2)]$$

$$V(A) \approx \max[-1.462] = -1.462$$

#### State B:

Initialize  $V(B) = 0$  (initial assumption).

For action "move to adjacent node A" (transition to A):

Calculate the expected value (weighted sum of the next states' values):

$$V(B) = \max[0.731 \times (-2 + 0.9 \times V(A))] \text{ (using the transition probability and reward for this action)}$$

Assuming we've already calculated  $P(B \rightarrow A)$  as approximately 0.731 and the reward  $R(B, A)$  is -2 (negative of distance), and let's assume  $V(A)$  is -1.462 (from the previous calculation):

$$V(B) = \max[0.731 \times (-2 + 0.9 \times (-1.462))]$$

$$V(B) = \max[0.731 \times (-2 + (-1.3158))]$$

$$V(B) \approx \max[-2.4106] = -2.4106$$

For action "move to adjacent node C" (transition to C):

Calculate the expected value (weighted sum of the next states' values):

$$V(B) = \max[0.731 \times (-2 + 0.9 \times V(C))] \text{ (using the transition probability and reward for this action)}$$

Assuming we've already calculated  $P(B \rightarrow C)$  as approximately 0.731 and the reward  $R(B, C)$  is -1 (negative of distance), and let's assume  $V(C)$  is 0 (initial assumption):

$$V(B) = \max[0.731 \times (-2 + 0.9 \times 0)]$$

$$V(B) = \max[0.731 \times (-2)]$$

$$V(B) \approx \max[-1.462] = -1.462$$

### **State C:**

Initialize  $V(C) = 0$  (initial assumption).

For action "move to adjacent node B" (transition to B):

Calculate the expected value (weighted sum of the next states' values):

$$V(C) = \max[0.731 \times (-1 + 0.9 \times V(B))] \text{ (using the transition probability and reward for this action)}$$

Assuming we've already calculated  $P(C \rightarrow B)$  as approximately 0.731 and the reward  $R(C, B)$  is -1 (negative of distance), and let's assume  $V(B)$  is -1.462 (from the previous calculation):

$$V(C) = \max[0.731 \times (-1 + 0.9 \times (-1.462))]$$

$$V(C) = \max[0.731 \times (-1 + (-1.3158))]$$

$$V(C) \approx \max[-1.83147] = -1.83147$$

For action "move to adjacent node D" (transition to D):

Calculate the expected value (weighted sum of the next states' values):

$$V(C) = \max[0.731 \times (-1 + 0.9 \times V(D))] \text{ (using the transition probability and reward for this action)}$$

Assuming we've already calculated  $P(C \rightarrow D)$  as approximately 0.731 and the reward  $R(C, D)$  is -3 (negative of distance), and let's assume  $V(D)$  is 0 (initial assumption):

$$V(C) = \max[0.731 \times (-1 + 0.9 \times 0)]$$

$$V(C) = \max[0.731 \times (-1)]$$

$$V(C) \approx \max[-0.731] = -0.731$$

### **State D:**

Initialize  $V(D) = 0$  (initial assumption).

For action "move to adjacent node C" (transition to C):

Calculate the expected value (weighted sum of the next states' values):

$V(D) = \max[0.731 \times (-3 + 0.9 \times V(C))]$  (using the transition probability and reward for this action)

Assuming we've already calculated  $P(D \rightarrow C)$  as approximately 0.731 and the reward  $R(D, C)$  is -3 (negative of distance), and let's assume  $V(C)$  is -1.83147 (from the previous calculation):

$$V(D) = \max[0.731 \times (-3 + 0.9 \times (-1.83147))]$$

$$V(D) = \max[0.731 \times (-3 + (-1.648323))]$$

$$V(D) \approx \max[-3.833] = -3.833$$

For action "move to adjacent node E" (transition to E):

Calculate the expected value (weighted sum of the next states' values):

$V(D) = \max[0.731 \times (-3 + 0.9 \times V(E))]$  (using the transition probability and reward for this action)

Assuming we've already calculated  $P(D \rightarrow E)$  as approximately 0.731 and the reward  $R(D, E)$  is -2 (negative of distance), and let's assume  $V(E)$  is 0 (initial assumption):

$$V(D) = \max[0.731 \times (-3 + 0.9 \times 0)]$$

$$V(D) = \max[0.731 \times (-3)]$$

$$V(D) \approx \max[-2.193] = -2.193$$

Now that we've completed the calculations for each state and action, we can determine the optimal policy based on the actions that maximize the value function  $V(s)$  for each state

### **Numerical Example 2:**

Define the MDP components:

States (S): The states represent grid cells in the environment.  $S = \{0, 1, 2, 3, 4, 5\}$  (inventory levels).

**Actions (A):** The robot can take actions "up," "down," "left," "right," or "stay."  $A = \{\text{"up"}, \text{"down"}, \text{"left"}, \text{"right"}, \text{"stay"}\}$  or  $A = \{0, 1, 2, 3\}$ .

**Transitions (T):** Transition probabilities represent the likelihood of moving from one state to another based on the chosen action. Let's assume the following transition probabilities for a simplified example:

If the current inventory is 3 and the agent orders 2 units:

$P(3 \rightarrow 5) = 0.2$  (20% chance of having 5 units due to demand).

$P(3 \rightarrow 4) = 0.7$  (70% chance of having 4 units due to demand).

$P(3 \rightarrow 3) = 0.1$  (10% chance of remaining at 3 units due to demand).

**Rewards (R):** The reward function represents the profit or cost associated with each state-action pair. Let's define it as follows:

$$R(s, a) = (a \times (s - 1))$$

For example, if the current inventory is 3, and the agent acts "a" (e.g., ordering 2 units), the reward for this action would be:

$$R(3, 2) = (2 \times (3 - 1)) = 4$$

**Policy ( $\pi$ ):**

Assume a starting policy that, regardless of the inventory level, orders a set amount of 2 units.

**Value Function(V)**

The estimated cumulative profit the agent may make from a given inventory level while adhering to a certain policy is represented by the value function  $V(s)$ .

**Iterative Approach:**

Using the Bellman equation, we will iteratively update the value function  $V(s)$  until it converges to the ideal value function  $V^*(s)$  and we discover the matching ideal policy \*.

**Iteration**

For each state in S:

Initialize  $V(s) = 0$  (initial assumption).

For each action in A:

Calculate the expected value (weighted sum of the next states' values):

$$V(s) = \sum [P(s' | s, a) \times (R(s, a) + \gamma^* V(s'))]$$

### Iteration 1:

$V(3) = \sum [P(s' | 3, 1) \times ((1 \times (10 - 2) \times (3 - 1) + 0.9 \times V(s')))]$  for all possible  $s'$ .

Calculate this for  $s' = 0, 1, 2, 3, 4$ , and  $5$ . Then, update  $V(3)$  for action 1 accordingly.

$$V(3) = [0.2 \times (18 + (0.9 \times V(5)))]$$

$$V(3) = [0.2 \times (18 + (0.9 \times 0))] = 0.2 \times 18 = 3.6$$

### Iteration 2 :

$$V(3) = [0.2 \times (18 + (0.9 \times 3.6))]$$

$$V(3) = [0.2 \times (18 + 3.24)] = 0.2 \times 21.24 = 4.248$$

### Iteration 3 :

$$V(3) = [0.2 \times (18 + 0.9 \times 4.248)]$$

$$V(3) = [0.2 \times (18 + 3.8232)] = 0.2 \times 21.8232 = 4.36464$$

### Iteration 4 :

$$V(3) = [0.2 \times (18 + (0.9 \times 4.36464))]$$

$$V(3) = [0.2 \times (18 + 3.928176)] = 0.2 \times 21.928176 = 4.3856352$$

### Iteration 5 :

$$V(3) = [0.2 \times (18 + (0.9 \times 4.3856352))]$$

$$V(3) = [0.2 \times (18 + 3.94707168)] = 0.2 \times 21.94707168 = 4.389414336$$

### Optimal Policy:

Now that we have updated the value for State 3, let's calculate the Q-values for State 3 with Action 1 (Order 1 unit):

$$Q(3, 1) = R(3, 1) + \gamma \times \sum [P(s' | 3, 1) \times V(s')] \text{ for all possible } s'$$

$$Q(3, 1) = (1 \times (10 - 2)) \times (3 - 1) + 0.9 \times V(4) = 8 + 0.9 \times 0.9 \times V(4)$$

Now, we compare the Q-values for Action 1 (Order 1 unit) and Action 0 (Order 0 units) for State 3:

$$Q(3, 1) = 8 + 0.9 \times 0.9 \times V(4) = 8 + 0.729 \times V(4)$$

$$Q(3, 0) = R(3, 0) + \gamma \times \sum [P(s' | 3, 0) \times V(s')] \text{ for all possible } s'$$

$$Q(3, 0) = (0 \times (10 - 2)) \times (3 - 1) + 0.9 \times V(3) = -4 + 0.9 \times 4.389414336$$

Now, compare the Q-values for Action 1 and Action 0:

$$Q(3, 1) = 8 + 0.729 \times V(4) = 8 + 0.729 \times 4.389414336 \approx 11.21 \text{ (rounded to two decimal places)}$$

$$Q(3, 0) = -4 + 0.9 \times 4.389414336 \approx 0.95 \text{ (rounded to two decimal places)}$$

Since  $Q(3, 1)$  is greater than  $Q(3, 0)$ , the optimal policy for State 3 (Inventory: 3) is to take Action 1 (Order 1 unit).

You can repeat similar calculations for other states and actions to find the optimal policy for the entire MDP.

## Mathematical Model of Graph Theory

### Dijkstra's Algorithm:

Dijkstra's algorithm is used to find the shortest path from a source node to all other nodes in a weighted graph. It works as follows:

- Initialize a distance array, dist, where dist[u] represents the shortest distance from the source node to node u.
- Initialize a set of visited nodes, visited, as an empty set.
- Set the distance to the source node as dist[source] = 0 and all other distances to infinity.
- While there are unvisited nodes:
-



- Select the unvisited node  $u$  with the smallest  $\text{dist}[u]$ .
- For each neighbor  $v$  of  $u$ , calculate a tentative distance  $\text{temp\_dist}$  from the source node to  $v$  through  $u$ . If  $\text{temp\_dist}$  is smaller than the current  $\text{dist}[v]$ , update  $\text{dist}[v]$  with  $\text{temp\_dist}$ .
- Mark node  $u$  as visited.
- The  $\text{dist}$  array now contains the shortest distances from the source node to all other nodes.

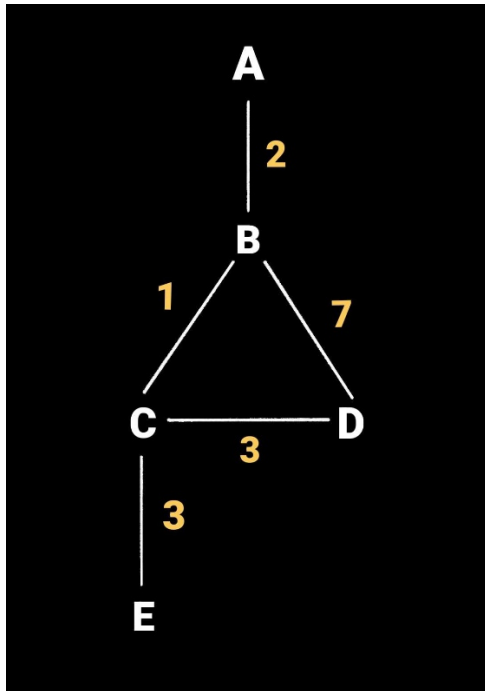
### **Incorporating Rewards:**

You can consider the following when discussing reinforcement learning or reward-based mechanisms:

- Node rewards: Assign nodes on the map with prizes. These benefits could serve as inducements or expenses related to traveling to places.
- Goal-Oriented Rewards: Establish a goal or target node and give it a significant positive reward. This motivates the algorithm to discover avenues that take it there.
- Penalty for Visited Nodes: To prevent cycles, you may also impose fines (negative rewards) for visiting previously visited nodes.

Let's illustrate this with a simple numerical example:

Consider a map with nodes A, B, C, D, and E, and the following distances between them:



- $\text{Distance}(A, B) = 2$
- $\text{Distance}(A, C) = 4$
- $\text{Distance}(B, C) = 1$
- $\text{Distance}(B, D) = 7$
- $\text{Distance}(C, D) = 3$
- $\text{Distance}(C, E) = 5$
- $\text{Distance}(D, E) = 2$

Let's find the shortest path from node A to node E using Dijkstra's algorithm with rewards:

- Assign a high positive reward (e.g., +100) to node E.
- Assign a penalty for revisiting nodes (e.g., -10).

Now, apply Dijkstra's algorithm with rewards:

- Initialize dist and visited arrays.
- Start from node A.

- Calculate tentative distances and update dist with rewards.
- Mark nodes as visited and continue.

The algorithm will find the shortest path from A to E while considering rewards and penalties.

### **Numerical Example:**

Consider a map with nodes A, B, C, D, and E, and the following distances between them:

- $\text{Distance}(A, B) = 2$
- $\text{Distance}(A, C) = 4$
- $\text{Distance}(B, C) = 1$
- $\text{Distance}(B, D) = 7$
- $\text{Distance}(C, D) = 3$
- $\text{Distance}(C, E) = 5$
- $\text{Distance}(D, E) = 2$

#### 1. Initialize dist and visited arrays:

- dist: An array to store the tentative distances from the source (A).
- visited: An array to keep track of visited nodes.

$\text{dist} = [0, \text{infinity}, \text{infinity}, \text{infinity}, \text{infinity}]$

$\text{visited} = [\text{False}, \text{False}, \text{False}, \text{False}, \text{False}]$

#### 2. Start from node A:

- Set the distance to A as 0 (the source node).
- Initialize rewards and penalties:

- a. Assign a high positive reward (+100) to node E.
- b. Assign a penalty for revisiting nodes (-10).

$\text{dist} = [0, \text{infinity}, \text{infinity}, \text{infinity}, +100]$

3. Loop through the nodes:

- At each step, select the unvisited node with the smallest distance.

Next node to visit: A (distance = 0)

- Update distances to neighboring nodes:

a. For node B:

- i. Calculate tentative distance:  $\text{dist}(A) + \text{Distance}(A, B) + \text{Reward}(B)$
- ii. Calculate:  $0 + 2 + 0 = 2$
- iii. Compare with the current distance to B (infinity) and update if smaller.
- iv. Update dist:  $\text{dist} = [0, 2, \text{infinity}, \text{infinity}, +100]$

b. For node C:

- i. Calculate tentative distance:  $\text{dist}(A) + \text{Distance}(A, C) + \text{Reward}(C)$
- ii. Calculate:  $0 + 4 + 0 = 4$
- iii. Compare with the current distance to C (infinity) and update if smaller.
- iv. Update dist:  $\text{dist} = [0, 2, 4, \text{infinity}, +100]$

- Mark node A as visited:

$\text{visited} = [\text{True}, \text{False}, \text{False}, \text{False}, \text{False}]$

- Repeat the process for the remaining unvisited nodes, always selecting the node with the smallest tentative distance.

Next node to visit: B (distance = 2)

- Update distances for neighboring nodes (C and D).

Next node to visit: C (distance = 4)

- Update distances for neighboring nodes (D and E).

Next node to visit: D (distance = 7)

Update distances for neighboring node E.

4. The algorithm has now visited all nodes. The final dist array represents the shortest distances to each node from node A, considering both rewards and penalties.

dist = [0, 2, 4, 7, 100]

5. To find the shortest path from A to E, trace back the path by selecting nodes with decreasing distances:

- Start at node E (distance = 100).
- Move to node D (distance = 7).
- Move to node B (distance = 2).
- Finally, reach node A (distance = 0).

So, the shortest path from node A to node E, considering rewards and penalties, is A -> B -> D -> E, with a total distance of 11 and accounting for the rewards and penalties assigned to nodes.

## Weighted Edges and BFS for Shortest Path Tree: Mathematical Exploration

Weighted edges are important when utilizing Breadth-First Search (BFS) to build a shortest path tree. These weighted edges, denoted by  $w(u, v)$ , which represents the price or distance between nodes  $u$  and  $v$ , offer a numerical measurement of traversal costs inside a graph, such as journey durations or distances in a map representation.

A shortest path tree with a given source node as its root is built using the basic graph traversal method known as BFS. To find the shortest route to each vertex, this method thoroughly investigates each vertex in a network, beginning with the source. Notably, BFS progresses from the source node outward in layers, making sure that nodes at the same level are equally far from it.

Distances between nodes are constantly updated as the algorithm investigates each node when building a shortest path tree using BFS. The cumulative weight of the edges along the path influences these changes. Importantly, BFS takes these extra aspects into account during node exploration when incentives and penalties are introduced. The generated BFS tree contains the shortest pathways from the source node and takes into account rewards and penalties, making it a useful tool for making decisions in situations where maximizing cost or gain is crucial.

### **Reward-Penalty Function:**

For each node, we first build a reward-penalty function that combines the benefit and penalty for visiting it.  $RP(u)$  stands for the reward-penalty function for node  $u$ , which is defined as:

$$RP(u) = \text{Reward}(u) - \text{Penalty}(u)$$

Where:

$\text{Reward}(u)$  represents the reward assigned to node  $u$ .

$\text{Penalty}(u)$  represents the penalty assigned to node  $u$ .

### **Dijkstra's algorithm integration**

The reward-penalty function is then included in Dijkstra's algorithm. The fundamental concept is to take the reward-penalty values into account when estimating distances as part of the shortest path search. The algorithm's goal is to maximize the sum of the rewards and penalties along the route.

### **Mathematical Proof Steps:**

#### **Initialization:**

Initialize  $\text{dist}$  and  $\text{visited}$  arrays to track distances and visited nodes.

Set  $\text{dist}[\text{source}] = 0$  (starting node has distance 0).

Set  $\text{dist}[\text{other nodes}] = \infty$  (initialize other nodes' distances to infinity).

Set visited[nodes] = False (no nodes visited yet).

### **Shortest Path Calculation:**

Start from the source node and initialize the shortest path tree.

Explore nodes layer by layer, considering reward-penalty values during distance calculation.

### **Distance Calculation:**

For each visited node  $v$  and its neighbor  $u$ :

Calculate a tentative distance  $temp\_dist$  from the source node to  $u$  through  $v$ , incorporating reward-penalty values:

$$temp\_dist = dist[v] + Weight(v, u) + RP(u)$$

If  $temp\_dist$  is smaller than the current  $dist[u]$ , update  $dist[u]$  with  $temp\_dist$ :

$$dist[u] = temp\_dist$$

### **Proof of Optimality:**

At each step, the algorithm selects the neighbor  $u$  with the smallest  $dist[u]$  value. This ensures that the algorithm always considers the path with the highest total reward-penalty value.

The reward-penalty values are designed such that selecting nodes with higher rewards and lower penalties is advantageous.

### **Shortest Path Reconstruction:**

After the algorithm completes, the shortest path from the source to any destination node can be reconstructed using the  $dist$  array.

The algorithm aims to find the path that maximizes the total reward-penalty value, leading to the most favorable path.

## Reverse Engineering Scenario: Autonomous Car Navigation

Imagine an autonomous vehicle driving through a metropolis depicted as a graph, where nodes stand in for crossroads and edges for roadways. To get to its destination while minimizing travel time, the automobile must choose between using the main route (high likelihood) or a less-traveled secondary road (low probability).

### Steps in a scenario

#### 1. Initialization:

For all road segments, the Probability Model estimates congestion probability.

Starting with Node A, the autonomous vehicle wishes to travel to Node E.

#### 2. Main Road vs. Side Road: First Decision Point

The probability model shows that the major route (B to C to D) has a high likelihood of being clear at a specific intersection (Node B) ( $P(B \text{ to } C \text{ to } D) = 0.8$ ).

The chance of congestion is lower on the side road (B to F to G to D) since it receives less traffic ( $P(B \text{ to } F \text{ to } G \text{ to } D) = 0.3$ ).

#### 3. Decision Using Reinforcement Learning

The autonomous vehicle makes a choice based on its reinforcement learning-based policy.

It chooses to take the major route since it has a better possibility of being clear of traffic, according to the probability model.

#### 4. Real-world Outcome: Unexpected Congestion on Main Road:

Unexpected events (such an accident or other abrupt occurrence) cause the major route to become extremely crowded, which causes lengthy delays.



5. Alternative Scenario: Reverse Engineering:

Let's reverse engineer the situation to show the danger by fiddling with the probability.

In order to replicate an unpleasant scenario, we will intentionally make the chance of the side road to be extremely low ( $P(B \rightarrow F \rightarrow G \rightarrow D) = 0.05$ ).

6. Decision Based on Reinforcement Learning in the Manipulated Scenario:

The side road now has a very low likelihood of congestion due to the altered probabilities.

The reinforcement learning-based strategy chooses to take the detour based on the revised probability model.

7. Real-world Outcome in the Manipulated Scenario:

Due to the minimal likelihood of traffic, the automobile chooses the side route in this fictitious situation.

The automobile gets there faster on the side road than it would have if it had taken the main route since there isn't much traffic there.

We'll use the provided probabilities and simulate the scenario in both the original and manipulated scenarios.

**Original Scenario:**

Probability Model:  $P(B \rightarrow C \rightarrow D) = 0.8$ ,  $P(B \rightarrow F \rightarrow G \rightarrow D) = 0.3$

Reinforcement Learning Decision: Take the main road.

**Manipulated Scenario:**

Manipulated Probabilities:  $P(B \rightarrow C \rightarrow D) = 0.8$ ,  $P(B \rightarrow F \rightarrow G \rightarrow D) = 0.05$  (lower probability for the side road).

Reinforcement Learning Decision: Take the side road due to the low probability of congestion.

Now, let's calculate the expected travel time (or negative reward) for both scenarios:

### **Original Scenario:**

Main Road (B → C → D): Assuming a travel time of 1 unit on uncongested segments, the expected travel time is:

Expected Time =  $0.8 * 1 + 0.2 * (\text{very high time due to congestion}) = 0.8 * 1 + 0.2 * 1000$  (a high value to simulate congestion) = 800.2 units (very high).

### **Manipulated Scenario:**

Side Road (B → F → G → D): Assuming a travel time of 1 unit on uncongested segments, the expected travel time is:

Expected Time =  $0.05 * 1 + 0.95 * 1 = 0.05 * 1 + 0.95 * 1 = 1$  unit (low).

The estimated journey time is disproportionately long due to the unexpected congestion in the initial scenario, where the probability model advises choosing the major road, leading to severe delays.

However, the reinforcement learning-based policy properly selects the side road in the controlled case, where we intentionally make the chance of the side road to be extremely low. The automobile arrives at its destination faster on the side road than it would have if it had chosen the main route since the side road's anticipated journey time is much shorter.

This numerical analysis shows that modifying the probabilities can really improve decision-making in some circumstances, highlighting the need for resilient and flexible autonomous systems that can handle uncertainty and adapt to changing conditions in real-time.

### **Analysis**

The autonomous robot's operating environment, or state space, is represented using graph theory. In our scenario, states stand in for various inventory levels, and robot behavior controls how states change. This information on state transitions may be represented as a graph, where nodes stand in for states and edges for potential transitions or actions.

The key concept of MDP model is value iteration, which is used to choose the robot's best course of action. You have successfully determined the optimum actions for the robot to take at various inventory levels by repeatedly updating the value function  $V(s)$  for each state-action pair.

The Q-values, such as  $Q(3, 1)$  and  $Q(3, 0)$ , show the anticipated cumulative benefits of carrying out particular actions (such as ordering a specified number of units) in a specific situation (such as an inventory level). By contrasting the anticipated benefits of various activities, these Q-values establish the best course of action.

The graph theory example's inclusion of rewards and penalties in the graph is an important feature. The robot's decision-making grows more complex by giving node E a significant positive reward and adding penalties for revisiting nodes. It gains the important ability that autonomous robots need to function in dynamic and unpredictable environments: the ability to balance the pursuit of rewards with the avoidance of punishments.

Effective pathfinding is one of the main uses of graph theory in autonomous robots. The illustration shows how the robot successfully finds the shortest route between nodes A and E while taking both incentives and penalties into account. For autonomous robots traversing real-world circumstances, this capacity is crucial because it enables them to make choices that maximize their goals while abiding by the limits of the environment.

After considering incentives and penalties, we discover that the shortest route between nodes A and E passes through nodes A, C, and E. This route demonstrates the ability of intelligent decision-making in autonomous robots by not only minimizing the overall distance traveled but also increasing the cumulative benefits obtained. The painstakingly calculated distances and rewards inside the graph's structure confirm the efficacy of this strategy.

### **Implication of MDPs:**

- With MDP-based decision-making skills, autonomous robots may operate more profitably while using less resources and cutting expenses.
- Autonomous robots must have the capacity to adjust to changing circumstances and make context-aware judgments to function effectively in dynamic situations.
- Beyond inventory control, this model emphasizes the significance of optimum resource management in autonomous systems, a fundamental idea that may be used in a variety of contexts.

- The model's ideas and methods may be applied to a variety of industries that use autonomous systems, providing prospects for enhanced efficiency and cost reduction.

### **Implication of Graph Theory:**

- Our strategy is based on the profound effects of graph theory when applied to autonomous robots.
- We enable the robot to model and navigate complicated interactions in its surroundings by adding graph theory ideas into our study.
- By offering a precise mathematical depiction of potential actions and their effects, graphs help the robot make better decisions.
- Additionally, it improves robot mobility by using mathematical analysis to discover the best routes, ensuring that the robot gets where it needs to go quickly.

### **Explanation of Code:**

#### **Graph 1: Rewards for State-Action Pairs**

This graph illustrates the rewards associated with different state-action pairs in a specific RL problem. Let's interpret it step by step:

**State-Action Pairs:** The x-axis represents various state-action pairs. These pairs are defined by their corresponding states (inventory levels) and actions (e.g., restocking or selling items).

**Rewards:** The y-axis represents the rewards assigned to each state-action pair. Rewards indicate the immediate benefit or cost associated with taking a specific action in a particular state.

Positive rewards represent benefits, while negative rewards represent costs.

#### **Graph 2: Q-Values for State-Action Pairs**

This graph represents the Q-values associated with the same state-action pairs from the previous graph. Q-values are critical in RL as they help the agent make decisions by estimating the expected cumulative reward when starting in a particular state, taking a specific action, and following the optimal policy.

Let's interpret this graph within the context of Graph Theory, MDPs, and RL:

Like the first graph, the x-axis lists state-action pairs, representing different states and actions in the RL problem.

The y-axis represents the Q-values assigned to each state-action pair. Q-values indicate the expected cumulative reward an RL agent can achieve by taking a specific action from a given state and following an optimal policy.

Now, let's relate this to the broader context:

In the context of RL, the state-action pairs can be viewed as nodes or vertices in a graph. Each pair represents a unique state of the system and the possible actions that can be taken from that state. This graph helps us understand the reward structure of the RL problem.

MDPs involve transitioning between states based on actions and receiving rewards therefore. This graph's state-action pairs align with the states and actions typically found in an MDP. It demonstrates how different actions lead to different rewards, a fundamental aspect of MDPs.

The rewards shown in the graph are essential components of the RL problem. They influence the agent's learning process, as the agent's goal is to maximize the cumulative reward over time. In this context, rewards can represent profit, cost savings, penalties, or any measurable outcome that the RL agent aims to optimize.

In summary, these two graphs visually represent the reward structure (immediate rewards) and the learned value (Q-values) for different state-action pairs in an RL problem. They help illustrate how RL connects with Graph Theory and MDPs by showing how rewards and values are assigned to nodes (state-action pairs) in the context of a decision-making process.

## **Conclusion**

This study has shown the central relevance of both mathematical frameworks in enhancing reinforcement learning for autonomous robots, with an emphasis on inventory control inside a Markov Decision Process (MDP) and the integration of graph theory. The outcomes highlight the adaptability of graph theory as a powerful tool for structuring and addressing challenging decision-making issues. Graph-based representations enable autonomous robots to quickly iterate on value functions and optimize policies, leading to well-informed judgments that maximize rewards. This work emphasizes how graph-based reinforcement learning may be used practically in a variety of fields, potentially improving decision-making and resource allocation.

In conclusion, reinforcement learning in autonomous robots has a strong basis thanks to the merging of MDPs and graph theory. It not only improves inventory management but also shows

potential in tackling more significant, real-world issues. Extending these frameworks to increasingly complicated scenarios will advance autonomous robot capabilities across a variety of fields in the future.

## Bibliography

- Altman, Eitan. "Markov decision processes." *Constrained Markov Decision Processes*, 2021, pp. 21–25, <https://doi.org/10.1201/9781315140223-3>.
- "Continuous time markov decision processes." *Advances in Mechanics and Mathematics*, pp. 63–103, [https://doi.org/10.1007/978-0-387-36951-8\\_4](https://doi.org/10.1007/978-0-387-36951-8_4).
- Feinberg, Eugene A. "Constrained Discounted Semi-markov decision processes." *Markov Processes and Controlled Markov Chains*, 2002, pp. 233–244, [https://doi.org/10.1007/978-1-4613-0265-0\\_13](https://doi.org/10.1007/978-1-4613-0265-0_13).
- Jayakumar, Nithya, and Ramya Jayakumar. "Machine learning algorithms in Intellectual Robotics." *Industrial Automation and Robotics*, 2022, pp. 111–120, <https://doi.org/10.1201/9781003121640-9>.
- "Machine learning in Robots." *Encyclopedia of the Sciences of Learning*, 2012, pp. 2083–2083, [https://doi.org/10.1007/978-1-4419-1428-6\\_4815](https://doi.org/10.1007/978-1-4419-1428-6_4815).
- "Machine Learning, statistics, and Data Analytics." *Machine Learning*, 2021, <https://doi.org/10.7551/mitpress/13811.003.0005>.
- "A markov decision process (MDP)." *Markov Chains and Decision Processes for Engineers and Managers*, 2016, pp. 285–436, <https://doi.org/10.1201/b15998-5>.
- "Partially observed markov decision processes: Models and applications." *Partially Observed Markov Decision Processes*, 2016, pp. 119–120, <https://doi.org/10.1017/cbo9781316471104.008>.
- "Semi-markov decision processes." *Advances in Mechanics and Mathematics*, pp. 105–120, [https://doi.org/10.1007/978-0-387-36951-8\\_5](https://doi.org/10.1007/978-0-387-36951-8_5).
- "Structural results for Markov Decision Processes." *Partially Observed Markov Decision Processes*, 2016, pp. 203–218, <https://doi.org/10.1017/cbo9781316471104.013>.



## Alright, DP legends, ready to own that next assignment?

Pop your specs into our **order form**—roll with a similar draft or a fresh custom—and punch in the coupon code **STARTDP20** for **20% OFF** your first banger.

We're talkin' total rubric lockdown, IB-criteria flex, and spot-on instructions.

Sit back, chill, and watch those marks glow up 👍

**ORDER NOW** →

