

# Analysis of Model-Free Reinforcement Learning Algorithm for Target Tracking

Muhammad Fikry<sup>1</sup>, Rizal Tjut Adek<sup>2</sup>, Zulfhazli<sup>3</sup>, Subhan Hartanto<sup>4</sup>, Taufiqurrahman<sup>5</sup>, Dyah Ika Rinawati<sup>6</sup>

<sup>1,2</sup>Department of Informatics, Universitas Malikussaleh, Indonesia

<sup>3</sup>Department of Civil Engineering, Universitas Malikussaleh, Indonesia

<sup>4</sup>Department of Computer Science, Universitas Pat Petulai, Indonesia

<sup>5</sup>Department of Informatics, Universitas Sumatera Utara, Indonesia

<sup>6</sup>Department of Industrial Engineering, Universitas Diponegoro, Indonesia

## ARTICLE INFO

### Article history's:

Received : 30 January 2022

Revised : 28 February 2022

Accepted : 17 March 2022

### Keywords:

Target tracking  
Machine learning  
Reinforcement learning  
Q-Learning

## ABSTRACT

Target tracking is a process that can find points in different domains. In tracking, some places contain prizes (positive or negative values) that the agent does not know at first. Therefore, the agent, which is a system, must learn to get the maximum value with various learning rates. Reinforcement learning is a machine learning technique in which agents learn through interaction with the environment using reward functions and probabilistic dynamics to allow agents to explore and learn about the environment through various iterations. Thus, for each action taken, the agent receives a reward from the environment, which determines positive or negative behavior. The agent's goal is to maximize the total reward received during the interaction. In this case, the agent will study three different modules, namely sidewalk, obstacle, and product, using the Q-learning algorithm. Each module will be training with various learning rates and rewards. Q-learning can work effectively with the highest final reward at a learning rate of 0.8 for 500 rounds with an epsilon of 0.9.

### Muhammad Fikry,

Department of Informatics, Universitas Malikussaleh, Indonesia

Email: muhammad.fikry@unimal.ac.id

## 1. INTRODUCTION

Determination of distribution channels can determine the optimal path with distance, time, and capacity as aspects of consideration [1]. Many methods can be used to optimize the path, such as the Heuristic Algorithm, the Tabu Search Method for delivery [2], the Ant Colony Optimization (ACO) method for mobile robot navigation [3], and others. Getting the optimal path will provide many advantages in terms of time, speed, cost and avoid losses caused by certain conditions such as congestion. The optimal path is determined from one place to a particular destination, and it has a starting point and stops at a predetermined endpoint with various conditions encountered. We try to approach the initial situations at any point and the stopping point in more than one place, meaning that it has several places as temporary stops and one final stop that occurs anywhere (random location).

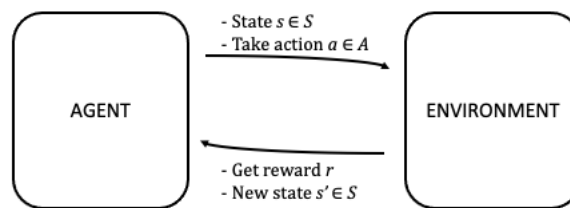
In this study, we consider the problem of tracking multiple objects in an environment. An environment is a rectangular space bounded by a sidewalk. We use model-free reinforcement learning to train the tracker to detect the product and avoid any obstacle in its movement. In other words, the agent must find the best way to achieve the highest reward. However, one of the challenges in reinforcement learning is the problem of exploration vs. exploitation, which is a trade-off between obtaining a tip from a choice that is considered safe and exploring other possibilities that may be profitable. The investigation allows increasing agent knowledge, which leads to long-term benefits. At the same time, exploitation chooses the action to get the most rewards but leads to less than optimal behavior. Our learning is done through python simulations to observe the system's characteristics. The agent will explore specific paths to get information. Then, the agent uses this information

in the exploitation to choose the best path for branching lanes. The agent will take action using greedy epsilon. Therefore, it is essential to get the best level of learning for agents to bring efficiency in taking steps.

## 2. METHOD

### 2.1. Reinforcement Learning

As part of machine learning methods, reinforcement learning (RL) [4,5] has been widely used in several disciplines to find optimal policies in an uncertain environment. Many studies have been carried out to develop RL techniques to find optimal feedback solutions, such as the zero-sum game problem [6] RL works by interacting with the environment. The RL algorithm works like the human brain when making several decisions and helps to take decisions sequentially. **Fig. 1** shows the interaction of agents with the environment in taking action to maximize reward.



**Fig. 1.** Interaction of agent with environment

RL studied effective strategies for agents from experimental trials and received feedback. With an optimal design, agents can actively adapt to the environment to maximize future rewards. The agent can remain in one of the many states ( $s \in S$ ) in the background and choose to take action ( $a \in A$ ) to move from one state to the next. The state where the agent will arrive is determined by the transition probability between states ( $P$ ). After the action is taken, the environment rewards ( $r \in R$ ) as feedback. In this study, state ( $s$ ) is the location, activity ( $a$ ) is the action (movement) that must be taken, and the reward ( $r$ ) is the initial value given when the move occurs.

In defining the reward function and transition probability. If the model is known, it means planning with perfect information, i.e., doing model-based RL. We can find the optimal solution with Dynamic Programming (DP) when we fully see the environment. However, if the model is unknown, that means learning with incomplete information, i.e., doing RL model-free or learning the model explicitly as part of the algorithm.

The model describes the environment, and we can study or infer how the environment will interact with models and provide feedback to agents. This model has two main parts, the transition probability function  $P$  and the reward function  $R$ . In equation (1), the transition function  $P$  records the transition probability from state ( $s$ ) to ( $s'$ ) after performing an action and obtaining ( $a$ ) temporary reward ( $r$ ).

$$P(s', r | s, a) = \mathbb{P}[S_{t+1} = s', R_{t+1} | S_t = s, A_t = a] \quad (1)$$

Thus, the transition function between states can be defined in equation (2) as a function of  $P(s', r | s, a)$ .

$$\begin{aligned} P_{ss'}^a &= P(s' | s, a) \\ &= \mathbb{P}[S_{t+1} = s', R_{t+1} | S_t = s, A_t = a] \\ &= \sum_{r \in R} P(s', r | s, a) \end{aligned} \quad (2)$$

And the reward function  $R$  in equation (3) predicts the next reward based on an action. The agent's only goal is to maximize the total reward he receives in the long run. Therefore, this reward function is important in mapping each state in the environment to a numerical reward.

$$\begin{aligned} R(s, a) &= \mathbb{E}[R_{t+1} | S_t = s, A_t = a] \\ &= \sum_{r \in R} r \sum_{r \in R} P(s', r | s, a) \end{aligned} \quad (3)$$

## 2.2. Q-Learning

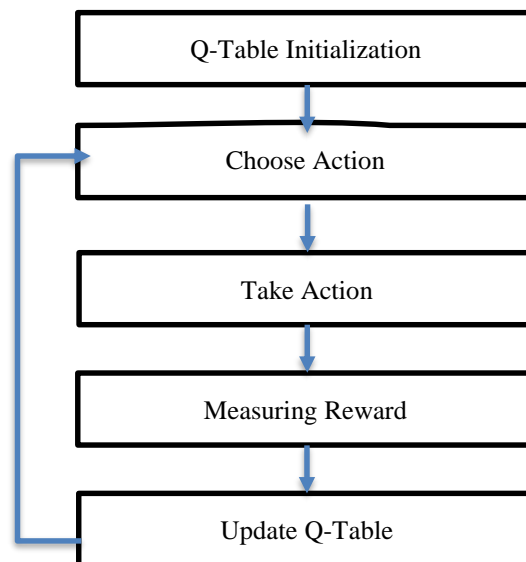
Q-learning is one of the important breakthroughs in reinforcement learning, which is the development of the Temporal Difference algorithm introduced by Watkins in 1989 [7,8] but The Q learning algorithm in its early use is flawed in several aspects and its application is limited [9]. Nowadays, some improved method of Q-Learning which combine with the Deep Learning method and named Deep Q-Learning and has been implemented in robotic application [10]. The Q in q-learning stands for quality, which indicates how useful a given action is in getting rewards in the future. Equation (4) below is in the form of the Q-learning formula.

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \max_{a \in A} Q(S_{t+1}, a) - Q(S_t, A_t)) \quad (4)$$

Q-learning takes action with the highest reward because Q-learning does not use a behavior policy to select the additional action  $A_{t+1}$ . Instead, estimate the expected future return in the update rule as  $\max A Q(S_{t+1}, A)$ . The max operator used follows a greedy policy. Q-Learning will converge to the optimal solution assuming that, after generating experience and training, it switches to a greedy policy. In equation (4),  $S_t$  is state  $S$  at time  $t$ ,  $A_t$  which is action  $A$  at time  $t$ , while  $S_{t+1}$  is the next state which will be the target of movement, and action  $A_{t+1}$  and reward  $R_{t+1}$ . Therefore, the transition probability function  $P$  and the reward function  $R$  affect the next reward that will be obtained. We define the set of states  $S$  as  $\{1, 2, \dots, 40\}$  which means we have 40 location points, the group of actions  $A$   $\{\text{up, down, left, right}\}$  which means that the possible moves are up, down, left, right, while for the positive  $R$  reward while  $\{1, 10, 100\}$  and the negative  $R$  reward while  $\{-1, -10, -100\}$ .

## 2.3. System Performance

The flow of system performance is shown in **Fig. 2** We will first create a table of  $n$  columns and  $m$  rows, where  $n$  is the number of actions and  $m$  is the number of states. Then assign zero values to all columns and rows. In the next step, we will select action  $a$  in state  $(s)$  based on the  $q$ -table, but as mentioned earlier, at the start, all values are 0. So, we use the greedy epsilon method to select actions randomly (up, down, left, right).



**Fig. 2.** System performance scenario using Q-learning

The system will execute the chosen action based on the greedy epsilon. After the move is obtained, the rewards and results are observed, then update the  $Q(s, a)$  function according to equation (4). After that, we will return to step two to do the next round until the condition stops. Table  $Q$  helps us find the best action in each state, and this is useful for maximizing the expected reward by choosing the best of all possible activities.  $Q(s,$

a) returns the expected future reward of the action on that state. When the q-table is ready, the agent will exploit the environment and take better steps.

### 3. RESULTS AND DISCUSSION

**Fig. 3** shows that when the agent moves, the agent will learn the best path to the end of the sidewalk through the reward function. Agent (A) starts at position (1,1). Sidewalks consist of rows and columns, not sidewalks denoted by a line (-). Obstacles (x) and Product (o) are placed randomly.



**Fig. 3.** Location of sidewalk, obstacles and product points

In this case, the state is a collection of locations, and actions are moves that the agent can perform. The state-transition function is the probability of an agent moving from one state to another through several actions. As an algorithm for learning with incomplete information, q-learning is chosen in this study for the application of RL, which is used to train agents in the environment. When q-learning is applied, we create a matrix that follows the form [state, action] and assigns an initial value (q-value) of zero. We then update and store the q-value after each round. The matrix, called q-table, becomes a reference table for agents to choose the best action based on the q-value. The state values for each possible action are entered into the q-table. The most significant matter is obtained according to the action that will give the maximum value from the next step.

We will train agents on each module and combine them linearly. The selected action is determined by the greedy method; the agent will use the appropriate q-table to choose the action with the maximum expected value. However, the steps are not deterministic, so it is possible that the agent will not initially take the best move for a given state to encourage exploration. The sidewalk module is trained to give positive rewards to agents when they reach the sidewalk's end and give them negative rewards when they move out of the sidewalk. In the obstacles module, the agent has a window with 3×3 dimensions. If the object is in the window in a specific state, the agent will get punishment based on the reward function based on the distance between (0 – 1). With this, if the agent is outside the sidewalk and is in window obstacles, the negative value obtained is even greater. It may result in the agent starting training again from the beginning. The opposite happens if the agent is in the product window in a particular state, the agent will get a reward based on the distance-based reward function (0 – 1).

Modules are trained with various learning rates and rewards. Learning rate determines the extent to which information can override old details. **Fig. 4** shows the last reward for each module in each training round. **Fig. 5** is a heatmap of the q-table on the sidewalk module for the “up” action. The darker areas show a lower q-value, while the lighter areas show a higher q-value. We display a heatmap from the q-table with the three highest learning rates for each module. **Fig. 5** shows that agents are learning to stay away from the sidewalk.

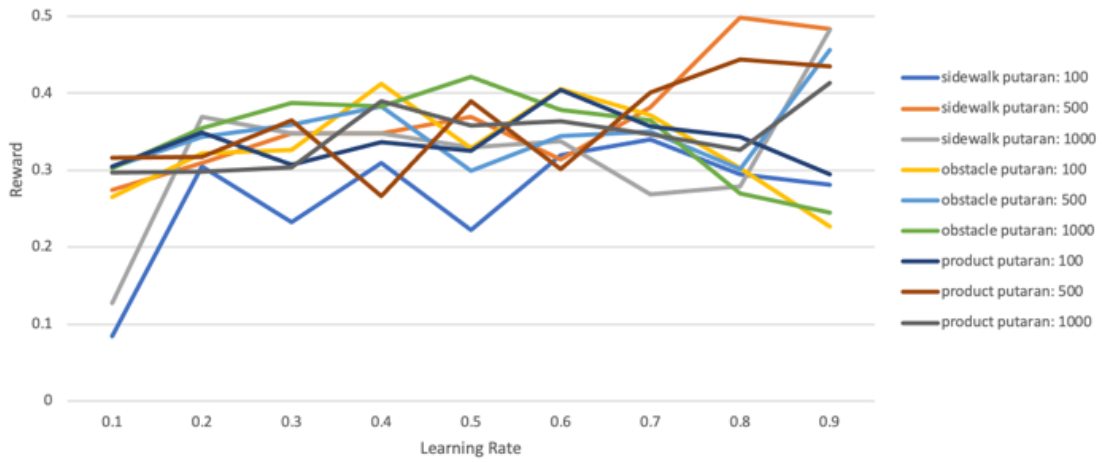


Fig. 4. The last reward for each module

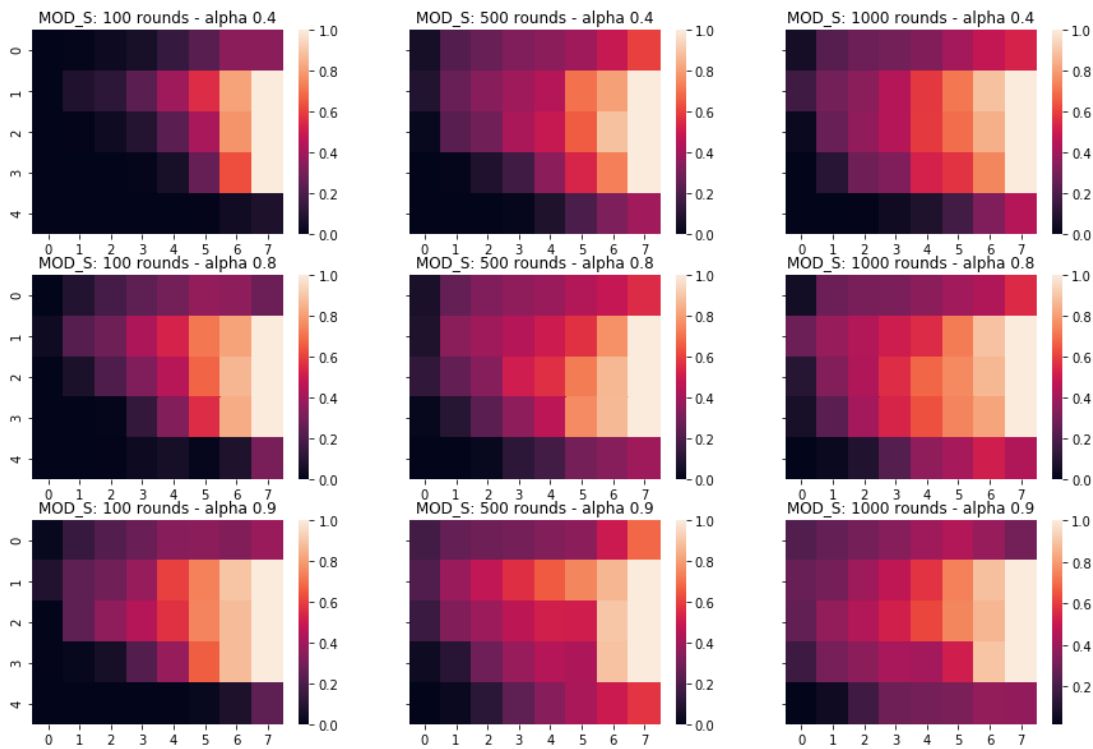


Fig. 5. Q-table on the sidewalk module

Furthermore, the obstacle module is trained with varying alpha values and rewards after 1000 training rounds. **Fig. 6** is a heatmap of the Q-table for the obstacle module for the “up” action. **Fig. 6** shows that the agent is learning to avoid obstacles.

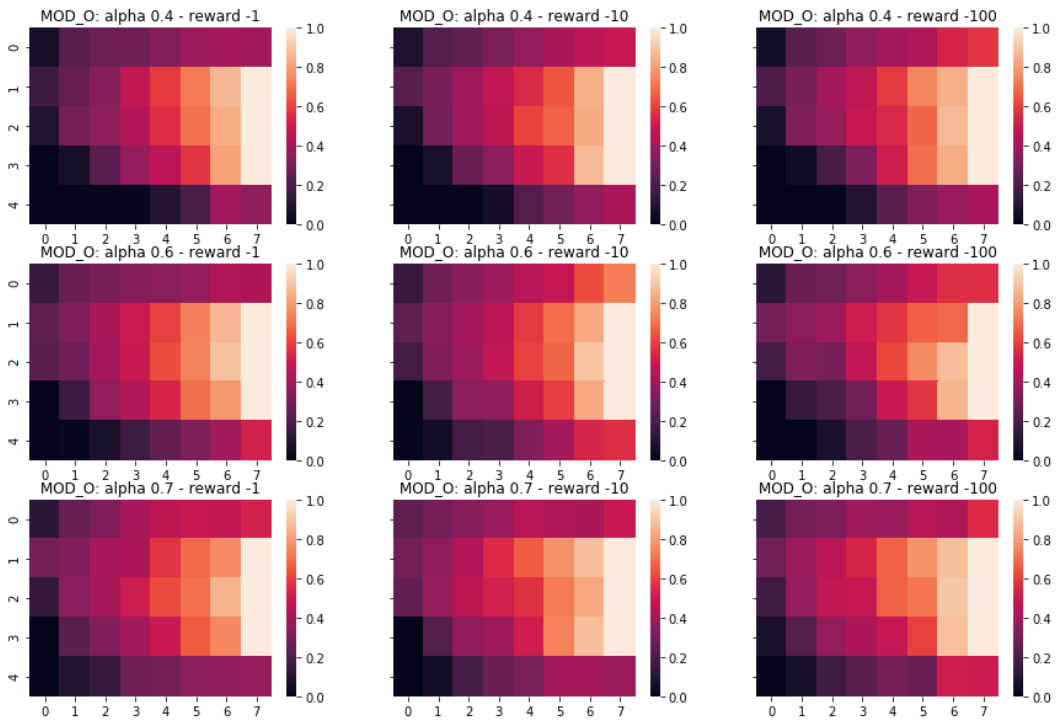


Fig. 6. Q-table on the obstacle module

Finally, the product module is trained on various alpha values and rewards after 1000 training rounds, and can be seen in Fig. 7.

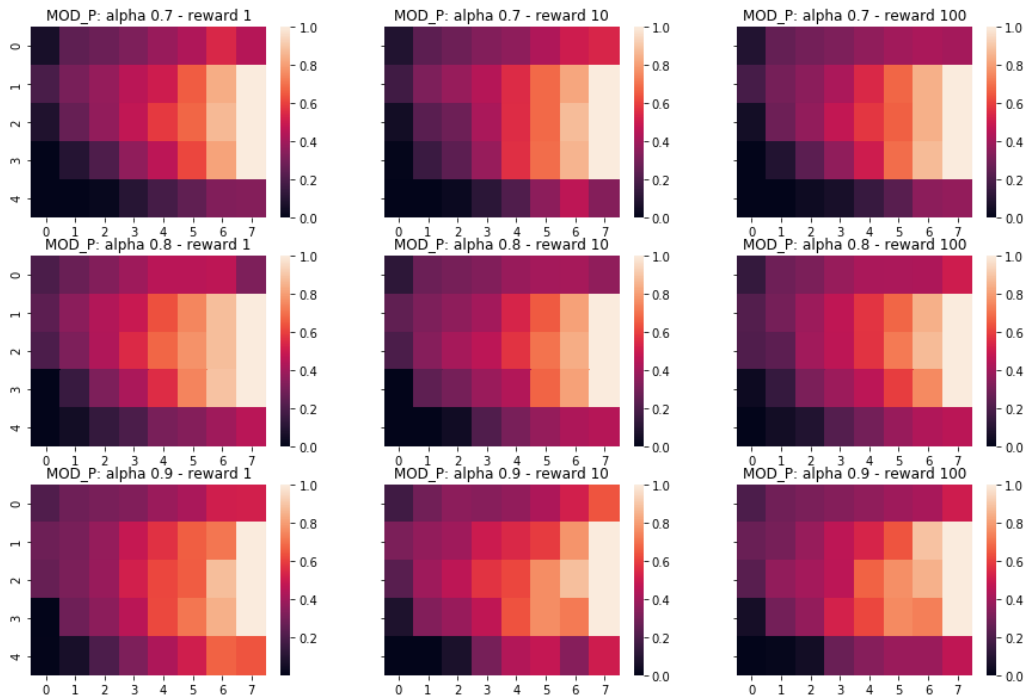
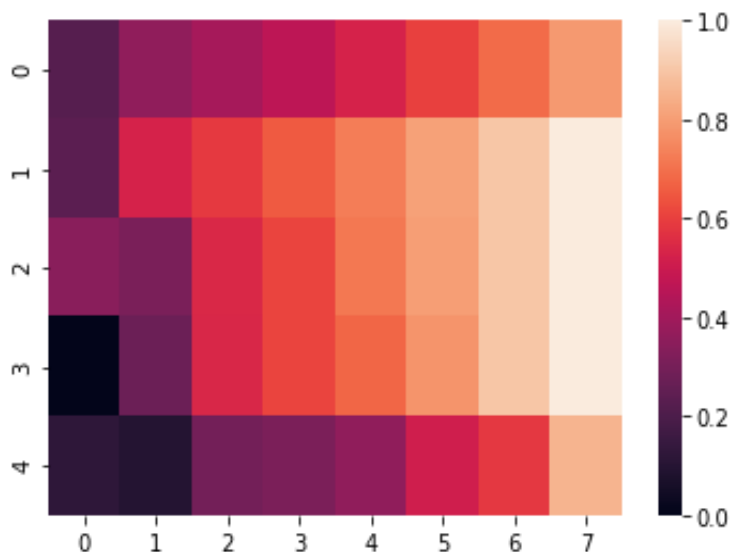


Fig.7. Q-table on the product module

**Fig. 8** is showing a heatmap of the linear combinations of the three q-tables for the “up” action.



**Fig 8.** Q-Table combination module

To balance exploration and exploitation during training, we used the epsilon-greedy method. We give a value of  $\epsilon=0.9$ , which means with a probability of 0.9, an action is chosen randomly from the action space. With a probability of 0.1, an effort is selected greedily based on  $\text{argmax}(Q)$ . The highest final reward is obtained with a learning rate of 0.8 in 500 rounds. More training provides more learning opportunities for agents, and this is the more stable the reward value will be. In other words, the agent is getting better at tracking, determining the action to the next state to get a positive reward.

#### 4. CONCLUSION

Q-learning is a practical algorithm for RL problems with incomplete information. This works best when the agent has a limited number of actions to perform in an environment with random conditions. Moreover, it effectively divides the problem into separate modules that can be adjusted and linearly combined. The agent will learn more to get the best reward with the higher learning rate value used. Still, the risk that arises is the possibility of error also increasing because it affects the agent to ignore previous knowledge to explore options. The higher the epsilon value, it allows the agent to expand his knowledge for each action, which leads to long-term benefits, increases the accuracy of value estimates, and allows the agent to make more informed decisions in the future, but may result in reduced chances of getting the most rewards. RL with the q-learning algorithm can also be used to determine the shortest route or optimize the movement of the robot in collecting rewards and avoiding certain states that can result in losses. The knowledge obtained from this research is expected to be a suggestion for system improvement or identification of influential variables related to tracking.

#### REFERENCES

- [1] Prabowo, R. C. (2018) ‘Penentuan Rute Distribusi Barang Yang Optimal Dengan Menggunakan Algoritma Heuristik pada PT. XYZ’, *Strategy: Jurnal Teknik Industri*, 3(1), pp. 47–50.
- [2] Hendra Bucika Glen Kadam, IG. Jaka Mulyana, J. M. (2018) ‘Penentuan Rute Terpendek Dengan Metode Tabu Search (Studi Kasus)’, *Widya Teknik*, 17(2), pp. 94–103.
- [3] Sirojul Hadi, Parama Diptya Widayaka, Radimas Putra M.D.L, R. D. (2020) ‘Pengukuran Jarak Pada Mobile Robot Menggunakan Xbee Berdasarkan Nilai Receive Signal Strength Indicator (RSSI)’, *Jurnal Bumigora Information Technology (BITE)*, 2(1), pp. 66–70.
- [4] Barto, R. S. S. and A. G. (2018) *Reinforcement Learning: An Introduction*. Second. MIT press.
- [5] Laber, J. C. And E. (2020) ‘Q-Learning: Theory and Applications’, *Annual Review of Statistics and Its Application*, 7, pp. 279–301.

- [6] Chakrabarti, S. M. O. T. S. (2018) ‘Reinforcement learning algorithms for uncertain, dynamic, zero-sum games’, in International Conference on Machine Learning and Applications (ICMLA). IEEE, pp. 48–54.
- [7] Watkins, C. J. C. H. (1989) Learning from delayed rewards. Cambridge University.
- [8] Ghulam Muhammad Ali, Asif Mansoor, Shuai Liu, Jacek Olearczyk, A. B. and M. A.-H. (2020) ‘Attaining global optimized solution by applying Q-learning’, in International Multidisciplinary Modeling & Simulation Multiconference, pp. 112–119.
- [9] B. Jang, M. Kim, G. Harerimana and J. W. Kim, "Q-Learning Algorithms: A Comprehensive Classification and Applications," IEEE Access, vol. 7, no. 1, pp. 133653-133667, 2019.
- [10] L. Tai and M. Liu, "A robot exploration strategy based on Q-learning network," in 2016 IEEE International Conference on Real-time Computing and Robotics (RCAR), Angkor Wat, 2016.
- [11] R.Liu, F. Nagotte, P.Zanne, M. Mathelin, and B. Dreps, “ Deep Reinforcement Learning fpr the Control of Robotic Manipulation: A Focussed Mini-Review”MDPI Journal of Robotics, 2021, 10, 22. <https://doi.org/10.3390/robotics10010022>.
- [12] T. Zhang and H. Mo. “Reinforcemnet Learning for Robot Research: A Comprehensive review and Open Issues” International Jurnal of Advanced Robotic Systems, 2021, 05, 01. DOI: 10.1177/17298814211007305.