# A Parameter Investigation of the $\epsilon$-greedy Exploration Ratio Adaptation Method in Multi-agent Reinforcement Learning

**Takuya Okano**                                              OKANO565656@GMAIL.COM

**Masaki Onishi**                                               ONISHI@NI.AIST.GO.JP

*Department of Policy and Planning Sciences*
*University of Tsukuba*
*Tsukuba, Ibaraki 305-0577, Japan*


*Artificial Intelligence Research Center*
*National Institute of Advanced Industrial Science and Technology*
*Tsukuba, Ibaraki 305-8560, Japan*


**Itsuki Noda**                                                 I.NODA@AIST.GO.JP

*Artificial Intelligence Research Center*
*National Institute of Advanced Industrial Science and Technology*
*Tsukuba, Ibaraki 305-8560, Japan*

## Abstract

Balancing the amounts of exploration and exploitation is one of the most important issues in reinforcement learning, which is also important while attempting to achieve favorable performances in multi-agent and non-stationary environments. In a previous study, we proposed the $\epsilon$-greedy exploration ratio $\epsilon$ adaptation method, which is referred to as "Win or Update Exploration ratio"(WoUE). This method is capable of adapting $\epsilon$ in multi-agent and non-stationary environments. However, while using this method, we have only one attempt to set the parameter that determines $\epsilon$'s update amount. In this study, we investigate the manner in which this parameter affects the WoUE's performance. Thus, we confirm that there is an optimal value for the WoUE's performance, and we further evaluate WoUE using the related adaptation method. We confirm that WoUE method improves the performance as compared with that obtained using the related methods regardless of the WoUE's parameter in multi-agent and non-stationary environments.

**Keywords:**   Reinforcement Learning, Multi-agent Reinforcement Learning, Exploration Ratio, $\epsilon$-greedy

## 1. Introduction

Exploration ratio adaptation in $\epsilon$-greedy is an important issue in reinforcement learning (Sutton and Barto, 1998). Several studies have proposed various adaptation methods, which only target single agent reinforcement learning and/or stationary environments. However, several difficulties arise while attempting to adapt these methods to scenarios with multi-agent and non-stationary environments. Therefore, we proposed exploration ratio adaptation method WoUE which adapts the exploration ratio in multi-agent and non-stationary environments. To function, this method requires a parameter that determines the exploration ratio update amount at one time. In this study, we investigate the influence that is

exhibited by various parameters on the system performance, and we evaluate WoUE and the related adaptation method of exploration ratio using each parameter.

## 2. Background

### 2.1 Domain

To evaluate the exploration ratio adaptation methods, we use the non-stationary "Repeated Resource Sharing Problem" (RRSP)(Noda and Ohta, 2008).

RRSP is a congestion game where agents share their resources. Mathematically, this RRSP can be represented by

$$
\begin{aligned}
RRSP &= \ <G, C, A, R>, \\
G &= \ \{g_1, g_2, \ldots, g_M\}, \\
C &= \ \{c_{g_1}, c_{g_2}, \ldots, c_{g_M}\}, \\
A &= \ \{a_1, a_2, \ldots, a_N\},
\end{aligned}
$$

where $G$ represents the set of resources, $C$ is the set of resource capacities, $A$ is the set of agents, $M$ is the number of resources, and $N$ is the number of agents. In addition, $n_j$, the number of agents that share the same resource $(j \in G)$, determines the reward function, $R(n_j, c_j)$, which is a monotonically decreasing function w.r.t. $n_j/c_j$. Each agent selects a resource $j$ according to $\epsilon$-greedy and gets a reward $R(n_j, c_j)$ by the selecting resource, and learns by reinforcement learning throughout the process, i.e., updates the expected reward $Q(j)$ by $Q(j) \leftarrow (1 - \alpha)Q(j) + \alpha R(n_j, c_j)$, where $\alpha$ is a learning ratio. This process is repeated during every step. The global optimal condition in the game is defined as the average of the distribution of all agent rewards is high value and the variance is minimum value. In other words, the global optimal condition is the capacity composition ratio of resource $j \in G$ and the number of agents, who select the resource $j$, is equal to a specific condition($\frac{c_j}{\sum_{k \in G} c_k} = \frac{n_j}{N} : \forall j \in G$). Therefore, we define the system performance, i.e., $Error$, using the following equation:

$$
Error \triangleq \sum_{j \in G} \left( \frac{c_j}{\sum_{k \in G} c_k} - \frac{n_j}{N} \right)^2. \tag{1}
$$

Here, the objective in the game is to minimize the $Error$ by multi-agent reinforcement learning as well as to evaluate the system performance using the $Error$.

We also introduce non-stationary environments to RRSP to reflect real-world situations. The capacity of each resource may change over time by a certain probability. We simulate numerous real-world problems using the non-stationary RRSP. On applying the congestion problem to RRSP, the road represents the resource and the car is an agent. We assume that the change in the resource capacity is an impassable or blocked road, consequently, RRSP creates a new road. Therefore, we use a non-stationary RRSP that resource's capacity change over time by a certain probability to evaluate the multi-agent reinforcement learning performance. Figure 1 shows the image diagram of the non-stationary RRSP.
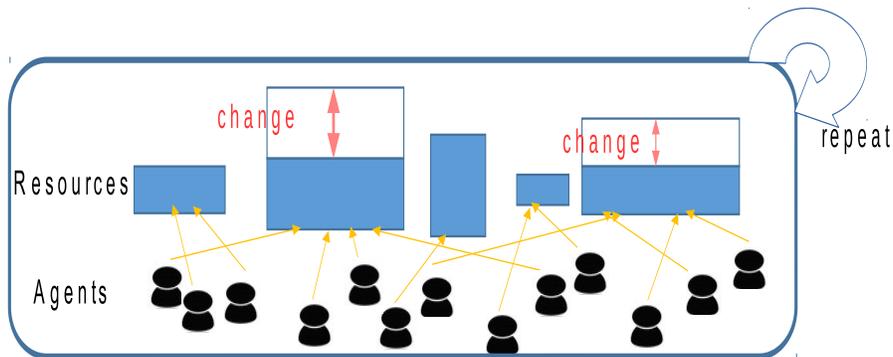
Figure 1: Diagram of the non-stationary RRSP.

## 2.2 Exploration Ratio Adaptation in Multi-agent and Non-stationary Environments

The exploration ratio $\epsilon$ in $\epsilon$-greedy is a parameter used to determine the balance between exploration and exploitation(Sutton and Barto, 1998). The ratio is a parameter that affects the learning performance as well as multi-agent reinforcement learning performance.

It is well known that this ratio has a positive value but should eventually converge to zero in a stationary environment. In non-stationary environments, however, this ratio should not converge to zero. Because the environment changes continuously over time, this ratio should be kept positive to follow the changes of the environment.

Moreover, in multi-agent environments, the agents influence each other. Therefore, the exploration ratio is also considered as a parameter that indicates the level of influence of an agent on the other agents. As a result, obtaining a suitable value for this parameter is important in order to improve performance of multi-agent reinforcement learning, which is a difficult yet essential issue.

Several previous studies have explored and proposed various methods to adjust the exploration ratio. However, these studies only targeted single agent reinforcement learning and/or stationary environments (Tokic, 2010; Tokic and Palm, 2011; Auer et al., 2002; Hester et al., 2013; Singh et al., 2000).

## 3. Win or Update Exploration Ratio

In our previous study, we proposed "Win or Update Exploration ratio"(WoUE) method (Okano and Noda, 2017), which is an $\epsilon$-greedy exploration ratio $\epsilon$ adaptation method in multi-agent and non-stationary environments.

In WoUE, agent $i$ who receives a lower reward compared with the average reward of all agents(total average) updates its $\epsilon_i$, with the objective of reaching equilibrium. The ratios are updated to minimize the variance of the distribution of all agents' rewards during a certain period according to the following equation:

$$\epsilon_i \leftarrow \begin{cases} \epsilon_i + \eta(\mu_R - R_i)(R_{\epsilon_i} - R_{(1-\epsilon_i)}) & (\mu_R > R_i) \\ \epsilon_i & (\mu_R \leq R_i), \end{cases} \tag{2}$$

3

where $\mu_R$ is the average reward for all agents respectively, $\epsilon_i$ is the exploration ratio of the agent $i$ for a certain period, $R_i$ is the average reward that the agent $i$ receives associated with the exploration ratio $\epsilon_i$ for a certain period. $R_{(1-\epsilon_i)}$ is the average reward derived from exploitation, and $R_{\epsilon_i}$ is the average reward derived from exploration. The update amount is given by the following equation:

$$(\mu_R - R_i)(R_{\epsilon_i} - R_{(1-\epsilon_i)}) \propto -\frac{\partial \sigma_R^2}{\partial \epsilon_i}. \tag{3}$$

Therefore, we observe that agent $i$ updates its exploration ratio, $\epsilon_i$, in an attempt to minimize the variance of the distribution of all agents' rewards $\sigma_R^2$. In other words, the agent $i$ updates $\epsilon_i$ to reach equilibrium. Because only the agents whose rewards are lower than total average, it can be expected to avoid falling into a sub-optimal equilibrium where the average reward is low. The derivation of Equation (3) is shown in Appendix A. Since each agent only requires $\mu_R$, this method is capable of updating $\epsilon$ in such a manner that minimizes $\sigma_R^2$ by only broadcasting $\mu_R$ to all agents. This method is easily adapted if there are many agents in the environment.

In this study, we use a different version of WoUE (Okano and Noda, 2017), which updates $\epsilon$ at each step in order to increase the opportunity to update $\epsilon$. To update at each step, the parameters $\mu_R$, $R_i$, $R_{(1-\epsilon_i)}$, and $R_{\epsilon_i}$ are modified with Exponential Moving Average. WoUE algorithm is shown in Appendix B.

## 4. Investigation of WoUE's Parameter $\eta$

WoUE initializes the parameter $\eta$ that modifies $\epsilon$ from Equation (2). We investigated the effects of $\eta$ on the $Error$(i.e., the system performance) via experimentation. The results indicate that there is an optimal $\eta$. When $\eta$ is high, WoUE can achieve a good performance.

We investigate changes of the $Error$ for each $\eta$ in the non-stationary RRSP. In this experiment, we compare the $Error$ with each $\eta$. Figure 2 shows the $Error$ of WoUE for
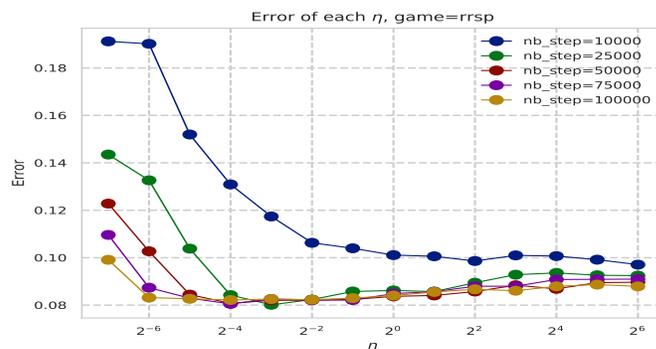


Figure 2: $Error$(i.e., the system performance) of WoUE associated with each $\eta$.

each $\eta$. Each line represents a different step number (the update time of the exploration ratio $\epsilon$). The results indicate that WoUE $Error$ curves over each $\eta$ have convex downward shapes, which confirms that there is an optimal $\eta$ parameter which minimizes the $Error$. When both $\eta$ and the step number are small, the system performance is poor. Figure 3 shows the changes in $\epsilon$ with respect to WoUE with a small(left) and a large(right) value of $\eta$. Figure 3 confirms that updating $\epsilon$ is quite slow. Poor performance due to the insufficient
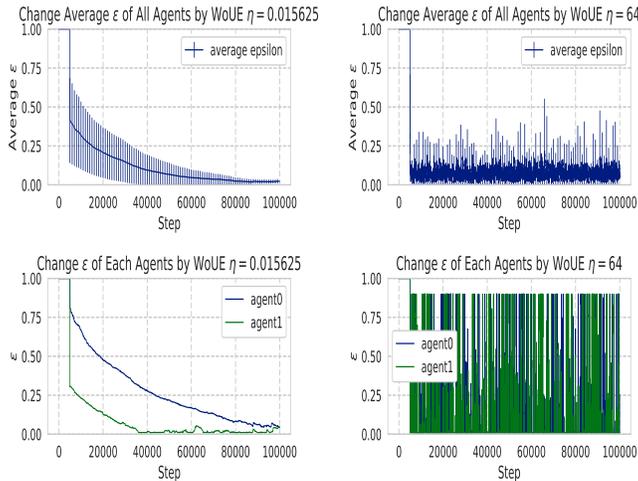
4

Figure 3: Change in $\epsilon$ by WoUE with $\eta = 2^{-6}$ (left) and $\eta = 2^6$ (right).The upper graphs are the average $\epsilon$ for all agents and the lower graphs are $\epsilon$ for the two sampled agents.

amount of time taken in this experiment to update $\epsilon$. If there is a sufficient amount of the update time, WoUE performs well. Although this method produces small *Error* quickly, when $\eta$ is large, the *Error* is bigger than the minimum value. Figure 3 indicates that $\epsilon$ changes rapidly, due to the large value of $\eta$, ($\epsilon$ transforms from a minimum $\epsilon$ to a maximum $\epsilon$, or from a maximum to a minimum). In addition, both the exploration noise and the *Error* increase.

Thus, if there is a sufficient amount of time to adapt $\epsilon$, $\eta$ should be set to a small value. When there is insufficient amount of time, $\eta$ should be set to a large value.

## 5. Comparison with Related Work

We compared WoUE with other methods in a non-stationary RRSP scenario. Based on these comparisons, we conformed that WoUE achieve a better performance than the other methods regardless of the various $\eta$ values in non-stationary environments.

### 5.1 Comparative Method

As comparative methods, we chose a normal $\epsilon$-greedy and "Value-Difference Based Exploration"(VDBE) method (Tokic, 2010). VDBE modifies $\epsilon$ in $\epsilon$-greedy method based on temporal-difference errors. Even though VDBE is used as an adaptation method for $\epsilon$ in single agent environments, this method is easily adapted to multi-agent environments. VDBE has two parameters $\sigma$ and $\delta$. Based on the results from (Tokic, 2010), we set $\delta = \frac{1}{M}$. In this experiment, we use a variety of $\sigma$ values and evaluate the normal $\epsilon$-greedy method without modifying $\epsilon$.

## 5.2 Experiment

We evaluate these methods using a non-stationary RRSP. This experiment comprises several agents ($N = 100$), many resources ($M = 5$), the reward function ($R(n_j, c_j) = \frac{c_j}{1+n_j}$), and the original resource capacities ($C = \{1, 2, 4, 8, 16\}$). Each capacity of resource doubles when a capacity is original and it restores to the original value when a capacity is already doubled by probability of 0.001 within each step. The parameters of each method are described in Appendix C. The learning ratio $\alpha$ for all agents is $\alpha = 0.3$. In these experiments, all agents act randomly for 5% of the steps such that all agents learn the environment.

| Method and Parameter | Average $Error$(SD) | Average $\epsilon$(SD) |
|---|---|---|
| $\epsilon = 0.01$ | 0.104(0.035) | 0.010(0.000) |
| $\epsilon = 0.1$ | 0.114(0.019) | 0.100(0.000) |
| VDBE $\sigma = 0.001$ | 0.569(0.069) | 0.788(0.032) |
| VDBE $\sigma = 0.01$ | 0.160(0.024) | 0.246(0.058) |
| VDBE $\sigma = 0.1$ | 0.156(0.058) | 0.001(0.003) |
| VDBE $\sigma = 1$ | 0.233(0.079) | 0.000(0.000) |
| WoUE $\eta = 0.015625$ | 0.109(0.021) | 0.090(0.030) |
| WoUE $\eta = 1$ | **0.095(0.023)** | 0.029(0.009) |
| WoUE $\eta = 64$ | 0.096(0.023) | 0.046(0.028) |

Table 1: Average $Error$ and average $\epsilon$ of all agents for each method in the non-stationary RRSP. SD is the standard deviation for $Error$ and $\epsilon$.
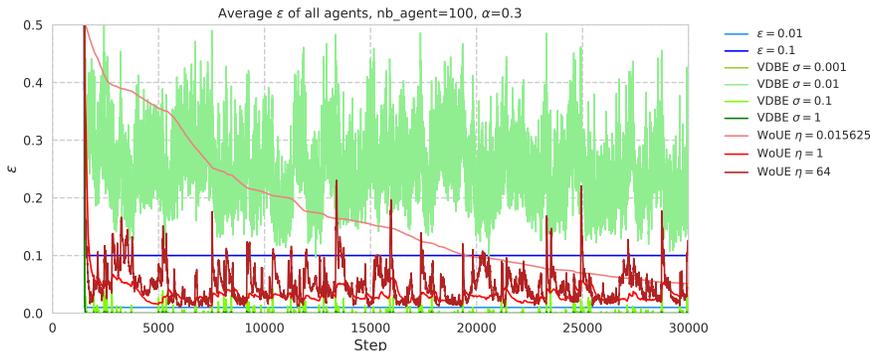


Figure 4: Changes in average $\epsilon$ of all agents for each method in non-stationary RRSP. VDBE with $\sigma = 0.001$ is out of this figure because the $\epsilon$ is very high.

Table 1 shows the $Error$(i.e., the system performance) and the average $\epsilon$ for all agents for each adaptation method in the latter half of all steps. The VDBE's $Error$ varies for each $\sigma$ value and the $Error$ is worse compared with the other methods, i.e., the VDBE's lowest $Error$ is higher than the WoUE's highest $Error$. On the other hand, the results indicate that WoUE has better performance than the other methods. Even though the VDBE's performance varies for each parameter, the WoUE's performance is favorable regardless of the parameters.

6

Figure 4 shows the change in the average $\epsilon$ for all agents. Here, the X-axis is the step number, and the Y-axis is the average $\epsilon$ for all agents. $\epsilon$ adapted by VDBE varies according to a specific parameter. This parameter affects the system performance sensitivity. Even though the SD of $\epsilon$ adapted by WoUE varies, the average ratio is similar. In another parameters, we confirmed the same result in Appendix D.

If WoUE's $\eta$ is small, the system performance is lacking due to a lowly changing $\epsilon$. However if we continually update $\epsilon$, system performance will improve.

## 6. Conclusion and Future work

In this study, we investigated the parameter $\eta$ in WoUE which is adaptation method of $\epsilon$ in $\epsilon$-greedy and compared the results with recent studies in multi-agent and non-stationary environments. The WoUE's parameter $\eta$ determines $\epsilon$'s the update amount. We conducted several experiments to investigate the effects of this parameter. Our results indicate that there are optimal $\eta$ values. We confirmed that when $\eta$ is small, WoUE achieves a better performance despite long simulation periods. When $\eta$ is high, even though WoUE achieves a better performance quicker, the system performance is worse than the best performance, due to the exploration noise. We evaluated WoUE with recent studies on VDBE with various parameters. Results indicate that WoUE performs more favorably than VDBE in multi-agent and non-stationary environments regardless of the parameters used.

The evaluation of WoUE in real world applications and its theoretical analysis are topics of future investigations.

## Acknowledgments

## Appendix A.

We further derive Equation(3). First, we express the average reward of the agent $i$ $R_i$ as

$$R_i \approx (1 - \epsilon_i)R_{(1-\epsilon_i)} + \epsilon_i R_{\epsilon_i}. \tag{4}$$

Equation(4) can be derived from the following set of equations:

$$
\begin{aligned}
R_i &= \frac{1}{T}\sum_{}^{T} r_t & (5)\\
&= \frac{1}{T}\left(\sum_{}^{T_{(1-\epsilon_i)}} r_{(1-\epsilon_i)} + \sum_{}^{T_{\epsilon_i}} r_{\epsilon_i}\right) & (6)\\
&= \frac{1}{T}\left(\frac{T_{(1-\epsilon_i)}}{T_{(1-\epsilon_i)}}\sum_{}^{T_{(1-\epsilon_i)}} r_{(1-\epsilon_i)} + \frac{T_{\epsilon_i}}{T_{\epsilon_i}}\sum_{}^{T_{\epsilon_i}} r_{\epsilon_i}\right) & (7)\\
&= \frac{T_{(1-\epsilon_i)}}{T}\frac{1}{T_{(1-\epsilon_i)}}\sum_{}^{T_{(1-\epsilon_i)}} r_{(1-\epsilon_i)} + \frac{T_{\epsilon_i}}{T}\frac{1}{T_{\epsilon_i}}\sum_{}^{T_{\epsilon_i}} r_{\epsilon_i} & (8)\\
&\approx (1-\epsilon_i)\frac{1}{T_{(1-\epsilon_i)}}\sum_{}^{T_{(1-\epsilon_i)}} r_{(1-\epsilon_i)} + \epsilon_i\frac{1}{T_{\epsilon_i}}\sum_{}^{T_{\epsilon_i}} r_{\epsilon_i} & (9)\\
&= (1-\epsilon_i)R_{(1-\epsilon_i)} + \epsilon_i R_{\epsilon_i}, & (10)
\end{aligned}
$$

$T_{1-\epsilon_i}$ represents the exploitation number and $T_{\epsilon_i}$ is the exploration number by agent $i$. $r_{(1-\epsilon_i)}$ is the reward from exploitation and $r_{\epsilon_i}$ is the reward from exploration.
Next we derive Equation(3). The variance $\sigma_R^2$ is calculated from:

$$\sigma_R^2 = \frac{1}{N}\sum_{i=0}^{N}(\mu_R - R_i)^2. \tag{11}$$

We differentiate $\sigma_R^2$ with respect to $\epsilon_i$

$$
\begin{aligned}
\frac{\partial \sigma_R^2}{\partial \epsilon_i} &\approx \frac{2}{N}(\mu_R - R_i)\frac{\partial(\mu_R - R_i)}{\partial \epsilon_i} & (12)\\
&\propto \frac{2}{N}(\mu_R - R_i)(R_{(1-\epsilon_i)} - R_{\epsilon_i}) & (13)\\
&\propto (\mu_R - R_i)(R_{(1-\epsilon_i)} - R_{\epsilon_i}) & (14)\\
-\frac{\partial \sigma_R^2}{\partial \epsilon_i} &\propto (\mu_R - R_i)(R_{\epsilon_i} - R_{(1-\epsilon_i)}). & (15)
\end{aligned}
$$

Therefore, we arrive at Equation(3):

$$(\mu_R - R_i)(R_{\epsilon_i} - R_{(1-\epsilon_i)}) \propto -\frac{\partial \sigma_R^2}{\partial \epsilon_i}. \tag{16}$$

## Appendix B.

---

**Algorithm 1** Multi-agent Reinforcement Learning with WoUE

---

1: Initialze N agents
2: **for** step $t = 0$ to T **do**
3:  **for** agent $i = 0$ to N **do**
4:   Select action, using $\epsilon$-greedy
5:  **end for**
6:  **for** agent $i = 0$ to N **do**
7:   Get reward $r_i$
8:   Update average reward $R_i \leftarrow \alpha r_i + (1 - \alpha)R_i$
9:   Update $Q_i(a)$
10:  **end for**
11:  Compute average reward of all agents $\mu_{R,t} = \frac{1}{N}\sum_{i=0}^{N} r_i$
12:  Update average reward of all agents $\mu_R \leftarrow \alpha\mu_{R,t} + (1 - \alpha)\mu_R$
13:  Broadcast $\mu_R$ to all agents
14:  **for** agent $i = 0$ to N **do**
15:   **if** agent $i$ exploited **then**
16:    Update exploited reward $R_{(1-\epsilon_i)} \leftarrow \epsilon_i r_i + (1 - \epsilon_i)R_{(1-\epsilon_i)}$
17:   **else if** agent $i$ explored **then**
18:    Update explored reward $R_{\epsilon_i} \leftarrow (1 - \epsilon_i)r_i + \epsilon_i R_{\epsilon_i}$
19:   **end if**
20:   **if** $R_i < \mu_R$ **then**
21:    $\epsilon_i \leftarrow \epsilon_i + \eta(\mu_R - R_i)(R_{\epsilon_i} - R_{(1-\epsilon_i)})$
22:    $\epsilon_i \leftarrow \text{CHECKMINMAX}(\epsilon_i)$
23:   **end if**
24:  **end for**
25: **end for**
26:
27: **function** CHECKMINMAX($\epsilon$)
28:  **if** $\epsilon < minimum\_\epsilon$ **then**
29:   $\epsilon \leftarrow minimum\_\epsilon$
30:  **else if** $maximum\_\epsilon < \epsilon$ **then**
31:   $\epsilon \leftarrow maximum\_\epsilon$
32:  **end if**
33:  **return** $\epsilon$
34: **end function**

---

9

## Appendix C.

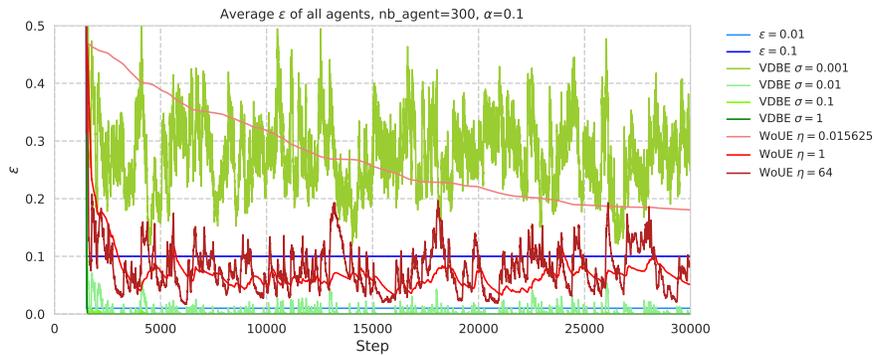| Name | Parameter details |
|------|-------------------|
| WoUE | $\eta = \{0.0125, 1, 245\}$, $minimum\_\epsilon = 0.01$, $maximum\_\epsilon = 0.9$ |
| VDBE | $\sigma = \{0.001, 0.01, 0.1, 1\}$, $\delta = \frac{1}{M}$ |

Table 2: Adaptation Methods for $\epsilon$ and their parameters.

## Appendix D.

Herein, we evaluated on the adaptation methods with other settings. This experiment comprises a number of agents ($N = 300$), a number of resources ($M = 8$), and the original resource capacities ($C = \{3, 6, 9, 12, 15, 18, 21, 24\}$). The learning ratio for all agents is $\alpha = 0.1$. Results are shown in  Table 3  and  Figure 5 . We confirmed that these results are similar to the result in Section 5.

| Method and Parameter | Average $Error$(SD) | Average $\epsilon$(SD) |
|----------------------|---------------------|------------------------|
| $\epsilon = 0.01$ | 0.112(0.057) | 0.010(0.000) |
| $\epsilon = 0.1$ | 0.091(0.034) | 0.100(0.000) |
| VDBE $\sigma = 0.001$ | 0.126(0.025) | 0.280(0.056) |
| VDBE $\sigma = 0.01$ | 0.170(0.067) | 0.004(0.006) |
| VDBE $\sigma = 0.1$ | 0.247(0.041) | 0.000(0.000) |
| VDBE $\sigma = 1$ | 0.300(0.051) | 0.000(0.000) |
| WoUE $\eta = 0.015625$ | 0.104(0.024) | 0.207(0.021) |
| WoUE $\eta = 1$ | 0.075(0.030) | 0.063(0.016) |
| WoUE $\eta = 64$ | **0.071**(0.032) | 0.077(0.037) |

Table 3: Average $Error$ and average $\epsilon$ of all agents for each method in the non-stationary RRSP. SD is the standard deviation for $Error$ and $\epsilon$.



Figure 5: Changes in average $\epsilon$ of all agents for each method in non-stationary RRSP. VDBE with $\sigma = 0.001$ is out of this graph.

## References

Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2–3):235–256, 2002.

Todd Hester, Manuel Lopes, and Peter Stone. Learning exploration strategies in model-based reinforcement learning. In *The Twelfth International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, May 2013.

Itsuki Noda and Masayuki Ohta. Meta-level control of multiagent learning in dynamic repeated resource sharing problems. In Tu-Bao Ho and Zhi-Hua Zhou, editors, *Proc. of PRICAI 2008*, pages 296–308. Springer, Dec. 2008.

Takuya Okano and Itsuki Noda. Adaptation method of the exploration ratio based on the orientation of equilibrium in multi-agent reinforcement learning under non-stationary environments. *Journal of Advanced Computational Intelligence and Intelligent Informatics*, 21(5):939–947, Sep. 2017.

Satinder Singh, Tommi Jaakkola, Michael L. Littman, and Csaba Szepesvári. Convergence results for single-step on-policy reinforcement-learning algorithms. *Machine Learning*, 38 (3):287–308, Mar 2000.

Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning : An Introduction*. MIT Press, 1998.

Michel Tokic. Adaptive $\epsilon$-greedy exploration in reinforcement learning based on value differences. In *Proceedings of the 33rd Annual German Conference on Advances in Artificial Intelligence*, KI'10, pages 203–210, Berlin, Heidelberg, 2010.

Michel Tokic and Gnther Palm. Value-difference based exploration: Adaptive control between epsilon-greedy and softmax. In *KI 2011: Advances in Artificial Intelligence*, volume 7006 of *Lecture Notes in Computer Science*, pages 335–346. Springer, 2011.