# 1 Hyper tuning using Q_learning

In order to tune the corresponding weights of distance, load , energy consumption, a variant of RL called Q learning is employed.

**Q-Learning principles according to the MEC system**:

1. **Environment:** The environment in which the agent operates is constructed from the set $[distance\_weight(W_d), load\_weight(W_l), EConsumption\_weight(W_{ec})]$ in range between 0 to 1 with step value of 0.05. As a consequence, it is represented as 3 dimensional matrix of size $n^3$ with n=20.

2. **Reward:** The $R_t$ rewards for each state in the environment is a metric that indicates the optimal weight assignment, which means that the particular combination of distance, load and energy consumption weights leads to lower application latency. It is constantly calculated from the average latency of mobile users(going through migration process).

$$R_t = 1/A_t$$

where, $A_t$ represents the average latency of mobile users till that timestamp since beginning. The above equation only show that lower the average latency, higher the reward value assign to the current combination values of $(W_d, W_l, W_{ec})$.

---

**Algorithm 1** Environment

---

1: **INPUT: Range** where, range of values between [0-1] (a multiple of 0.05)
2: **OUTPUT:** Rewards     /* Reward Matrix */
3: **for** $W_d$ $\epsilon$ **Range do**
4:     **for** $W_l$ $\epsilon$ **Range do**
5:         **for** $W_{ec}$ $\epsilon$ **Range do**
6:             **if** $W_d + W_l + W_{ec} = 1$ **then**
7:                 **if** current calculated $A_t$ = Ideal $A_t$ **then**
8:                     /* yet to decide the value of MAX($R_t$) */
9:                     rewards$[W_d][W_l][W_{ec}] \leftarrow$ MAX($R_t$)
10:                 **end if**
11:                 rewards$[W_d][W_l][W_{ec}] \leftarrow 1/A_t$
12:             **end if**
13:         **end for**
14:     **end for**
15: **end for**
16: **return** rewards

---

**The training of the agent on the environment**:

The training phase allows the agent to exploit the environment through the different actions it takes, by trying to improve its action strategy according to the rewards received.

1. **Action:** From any state $W_d, W_l, W_{ec}$, the agent has the possibility to execute the following 6 actions: stay in $W_d(|W_d)$, increment $W_d(+W_d)$, decrement $W_d(-W_d)$, stay in $W_l(|W_l)$, increment $W_l(+W_l)$, decrement $W_l(-W_l)$, stay in $W_{ec}(|W_{ec})$, increment $W_{ec}(+W_{ec})$, decrement $W_{ec}(-W_{ec})$(in multiple of 0.05). These actions are inserted in an Actions table, such as:

$$\textbf{Actions} = [' | W_d', ' + W_d', ' - W_d', ' | W_l', ' + W_l',$$
$$' - W_l', ' | W_{ec}', ' + W_{ec}', ' - W_{ec}]$$

2. **Q_table:** The Q_table contains the Q values for each action-state (Q(s, a)), such that the lines of the table are the states that represent the set of combinations $W_d, W_l, W_{ec}$ existing in the rewards matrix(algorithm 1), and the columns are the 6 actions that is described in Eq. 12. The Q values are initialized to 0 at the beginning of the training.

3. **Training functions:** To start the training, the agent must choose a random state $(W_d, W_l, W_{ec})$ in the environment (Algorithm 3). The state chosen by the agent must not be a terminal state that has a reward equal to $MAX(R_t)$ as indicated in algorithm 2. After choosing the departure state s $= W_d, W_l, W_{ec}$, the agent must perform an action according to the value of $\epsilon$ which indicates whether it is a random action or an action from the $Q\_table$ which maximizes the value Q(s, a) (algorithm 4). The execution of the action a at time t by the agent allows him to move to another state $s'$ and receive a reward $R_t$ according to the chosen action, as shown in algorithm 5.

4. **Agent Training:** Algorithm 6 gives the steps followed by the agent during its training with the parameters, fixed after multiple simulations, below, as well as the update of the $Q\_table$ at each state changes using the reward new state $s'$. This is done using the reward $R_t$ and the Q value of the current state Q($s'$, a); the agent calculates the time difference (TD) that will be used for the computation of the new Q value of the previous state s (Q(s,a)).

   - $\epsilon$: Exploration and exploitation time of environment(set as 0.9)
   - $\gamma$: Discount factor(set as 0.9)
   - $\alpha$: Learning Rate(set as 0.9)
   - Get_Next_Action($W_d$, $W_l$, $W_{ec}$, $\epsilon$): Allows to determine the action to be executed that has the largest Q value in the Q_table.
   - Get_Next_State($W_d$, $W_l$, $W_{ec}$, Action_Index): Allows to determine the next state from the executed action.

---
**Algorithm 2** Is_Termination_State
---
1: **INPUT**: $(W_d, W_l, W_{ec})$     // Test if the state is terminal
2: **OUTPUT**: Boolean Value
3: **if** rewards$[W_d][W_l][W_{ec}] = MAX(R_t)$ **then**
4:     return $\leftarrow$ True
5: **end if**
6: return $\leftarrow$ False

---

---

**Algorithm 3** Get_Starting_State

---

1: **INPUT**: range of values between [0-1].
2: **OUTPUT**: $W_d, W_l, W_{ec}$
3: For each $W_d, W_l, and\ W_{ec}$ choose a random real number(a multiple of 0.05) between 0 and 1;
4: **if** Is_Termination_State($W_d, W_l, W_{ec}$) **then**
5:     **return** Get_Starting_State()
6: **end if**
7: **return** $(W_d, W_l, W_{ec})$

---

---

**Algorithm 4** Get_Next_Action

---

1: **INPUT**: $(W_d, W_l, W_{ec}, \epsilon)$
2: **OUTPUT**: Action
3: **if** random() $< \epsilon$ **then**
4:     /* Choose the action that maximizes $Q_{value}$ for $(W_d, W_l, W_{ec})$ */
5:     **return** action with max($Q_{values}[W_d, W_l, W_{ec}]$)
6: **end if**
7: **return** random action from actions table

---

---

**Algorithm 5** Get_Next_State

---

1: /* current status and the selected action to perform */
2: **INPUT**: $(W_d, W_l, W_{ec}, Action\_Index)$
3: **OUTPUT**: $new\_W_d, new\_W_l, new\_W_{ec}$
4: **if** Actions[Action_Index]= $|W_d$ **then**
5:     $new\_W_d \leftarrow \bar{W}_d$
6: **end if**
7: **if** Actions[Action_Index]= $+W_d$ and $W_d <= 1$ **then**
8:     $new\_W_d \leftarrow \bar{W}_d + 0.05$
9: **end if**
10: **if** Actions[Action_Index]= $-W_d$ and $W_d >= 0$ **then**
11:     $new\_W_d \leftarrow \bar{W}_d$ - 0.05
12: **end if**
13: **if** Actions[Action_Index]= $|W_l$ **then**
14:     $new\_W_l \leftarrow \bar{W}_l$
15: **end if**
16: **if** Actions[Action_Index]= $+W_l$ and $W_l <= 1$ **then**
17:     $new\_W_l \leftarrow \bar{W}_l + 0.05$
18: **end if**
19: **if** Actions[Action_Index]= $-W_l$ and $W_l >= 0$ **then**
20:     $new\_W_l \leftarrow \bar{W}_l$ - 0.05
21: **end if**
22: **if** Actions[Action_Index]= $|W_{ec}$ **then**
23:     $new\_W_{ec} \leftarrow \bar{W}_{ec}$
24: **end if**
25: **if** Actions[Action_Index]= $+W_{ec}$ and $W_{ec} <= 1$ **then**
26:     $new\_W_{ec} \leftarrow \bar{W}_{ec} + 0.05$
27: **end if**
28: **if** Actions[Action_Index]= $-W_{ec}$ and $W_{ec} >= 0$ **then**
29:     $new\_W_{ec} \leftarrow \bar{W}_{ec}$ - 0.05
30: **end if**
31: **return** $new\_W_d, new\_W_l, new\_W_{ec}$

---

**Algorithm 6** Agent_Training

1: $W_d, W_l, W_{ec} \leftarrow$ Get_Starting_State()
2: **for iteration** $\rightarrow$ [1 - total number of migration] **do**
3:     **while** Is_Termination_State($W_d, W_l, W_{ec}$) = False **do**
4:         $Action\_Index \leftarrow$ Get_Next_Action($W_d, W_l, W_{ec}, \epsilon$)
5:         $new\_W_d, new\_W_l, new\_W_{ec} \leftarrow$ Get_Next_State($W_d, W_l, W_{ec}, Action\_Index$)
6:         reward $\leftarrow$ rewards[$new\_W_d$][$new\_W_l$][$new\_W_{ec}$]
7:         $old\_Q\_value \leftarrow Q\_values[W_d, W_l, W_{ec}, Action\_Index]$
8:         TD $\leftarrow$ reward + ($\gamma * MAX_a(Q\_values[W_d, W_l, W_{ec}])$) - $old\_Q\_value$
9:         $new\_Q\_value = old\_Q\_value + \alpha$*TD
10:        $Q\_values[W_d, W_l, W_{ec}, Action\_Index] \leftarrow new\_Q\_value$
11:     **end while**
12: **end for**