## Table 1: Decision Making

| Types | Value Function | Reward | Next Node [*] |
|---|---|---|---|
| sim | $Q- = \dfrac{expect - true}{expect}$ <br><br> $value\_exp = \dfrac{e^{Q_i}}{\sum_{j=1}^{k} e^{Q_j}}$ | – | $[expect \times$ <br> $value\_exp]$ |
| rl | $Q^{new}(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha.[r_t +$ <br> $\gamma.\max\limits_{a} Q(s_{t+1}, a) - Q(s_t, a_t)]^{**}$ <br> $value\_exp = \dfrac{e^{Q_i}}{\sum_{j=1}^{k} e^{Q_j}}$ | $r_t = log(|expect-$ <br> $true|)^{***}$ | $[expect \times$ <br> $value\_exp]$ |
| exp | $Q^{new}(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha.[r_t +$ <br> $\gamma.\max\limits_{a} Q(s_{t+1}, a) - Q(s_t, a_t)]$ <br> $value\_exp = \dfrac{e^{Q_i} * expect_i}{\sum_{j=1}^{k} e^{Q_j} * expect_j}$ | $r_t = log(|expect-$ <br> $true|)$ | $[value\_exp]$ |
| rl_true | $Q^{new}(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha.[r_t +$ <br> $\gamma.\max\limits_{a} Q(s_{t+1}, a) - Q(s_t, a_t)]$ <br> $value\_exp = \dfrac{e^{Q_i}}{\sum_{j=1}^{k} e^{Q_j}}$ | $r_t = true$ | $[expect+$ <br> $\max\limits_{\mathcal{N}}($ <br> $value\_exp)]^{****}$ |

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \Big( \underbrace{\overbrace{\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max\limits_{a} Q(s_{t+1}, a)}_{\text{estimate of optimal future value}}}^{\text{temporal difference}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \Big)$$
$$\underbrace{\phantom{r_t + \gamma \cdot \max Q(s_{t+1}, a) - Q(s_t, a_t)}}_{\text{new value (temporal difference target)}}$$

[*] Choosing neighbour with maximum value of $-$. Epsilon-greedy with $\epsilon = 0.1$
[**] For this simulation, $\alpha = 0.1, \gamma = 0.95$
[***] If $expect = true, r_t = 0$
[****] ( expected idleness + max($value\_exp$ of next neighbours) )


## Table 2: Estimating Idleness

| Types | Expression |
|---|---|
| sim | AGENT MODEL <br> Expected idleness = Time elapsed since last visit of that bot to that node |
| avg | OBSERVATION MODEL 1 <br> Expected Idleness = Average of all true idleness <br> observed on the previous visits <br> $expect_n = \dfrac{1}{n-1} * \sum_{j=1}^{n-1} true_j$ |
| ses [*] | OBSERVATION MODEL 2 <br> Expected Idleness = Simple Exponential Smoothing (SES) of all <br> true idleness observed on previous visits <br> $expect_n = \alpha\, true_{n-1} + \alpha(1-\alpha) true_{n-2} + \alpha(1-\alpha)^2 true_{n-3}...$ |

[*] $\alpha = 0.9$