
Algorithm 1 Pseudeocode of the Proposed Method

- 1: Establish a pair of DQNs (a trained DQN with weights θ and a target DQN with weights θ^-) and an empty experience pool M for all CUs;
 - 2: Initialize the trained DQN with random weights; set $\theta^- = \theta$;
 - 3: In time slot $t(M \leq M_b)$, CU takes action randomly and stores the corresponding experience $\langle s, a, r, s' \rangle$ in its experience pool.
 - 4: **repeat**
 - 5: CU k observe its state s_k in time slot t , $\forall k \in K$;
 - 6: In time slot $t(t > M_b)$ CU k chooses an action a_k according to ϵ -greedy policy; agent k chooses an action $a_k = \arg \max_{a \in A} q(s_k, a; \theta_k)$ with probability $(1 - \epsilon)$, or randomly chooses an action $a_k \in A$ with probability ϵ , $\forall k \in K$;
 - 7: CU k executes the chosen action a_k , then gets an immediate reward $r_k = R(s_k, a_k)$, $\forall k \in K$;
 - 8: CU k observes a new state s'_k in time slot $t + 1$, $\forall k \in K$;
 - 9: CU k saves its new experience $\langle s_k, a_k, r_k, s'_k \rangle$ into experience pool;
 - 10: Samples a mini-batch consisting of M_b experiences from experience pool M ;
 - 11: Updates the weights θ of trained DQN according to BP base on difference between reward in experience pool and output Q-value of target DQN;
 - 12: Updates θ^- with θ every T_{step} time slots;
 - 13: **until** $R(t) - R(t - 1) < \gamma$
-