**Algorithm 2:** Pseudocode: Deep Q Learning

---

1 Input: Initial $\epsilon$, Final $\epsilon$, $\epsilon$-decay rate, learning_rate $\alpha$, Reward Discounting Factor $\gamma$, Number of Steps $n\_steps$, $n\_rollouts$, $n\_iterations$, $batch\_size$, $n\_epochs$, $target\_update\_freq$, Environment

2 Initialize Replay Memory $\mathcal{RM}$ to capacity N

3 Initialize $Q$ and $\hat{Q}$ Network with random weights $w$

4 $total\_steps = 0$

5 $k = 0$

**Procedure** `rollout()`

1     Reset Environment and get $x_k$

2     k = 0

    **while** $k < n\_rollouts$ **do**

3         $u_k \leftarrow \begin{cases} \arg\max_u Q(x_k, u) & \text{probability } 1 - \epsilon \\ \text{Random action} & \text{probability } \epsilon \end{cases}$

4         Step Environment with action $u_k$ and get $x_{k+1}$ and $r_k$

5         Add transition $(x_k, u_k, r_k, x_{k+1})$ to $\mathcal{RM}$

6         $total\_steps \leftarrow total\_steps + 1$

        **if** $total\_steps \% target\_update\_freq == 0$ **then**

7             $\hat{Q} \leftarrow Q$

        **end**

8         $\epsilon \leftarrow \epsilon \times (decay)$

9         $k \leftarrow k + 1$

    **end**

    **return**

**Procedure** `learn()`

1     $i = 0$

    **while** $i < n\_epochs$ **do**

2         Sample random $batch\_size$ number of transitions$(x_j, u_j, r_j, x_{j+1})$ from $\mathcal{RM}$

3         $\hat{y}_j = \begin{cases} r_j & \text{if } x_{j+1} \text{ is terminal} \\ r_j + \gamma \max_u \hat{Q}_u(x_j, u) & \text{if } x_{j+1} \text{ is not terminal} \end{cases}$

4         $y_j = Q(x_j, u_j)$

5         Perform gradient descent on $(\hat{y}_j - y)^2$

6         $i \leftarrow i + 1$

    **end**

    **return**

    **while** $k < n\_iterations$ **do**

6     `rollout()`

7     `learn()`

8     $k \leftarrow k + 1$

    **end**