```
10: for m = 1 ... MAX_STEPS do
11:     Gather and store h experiences (s_i, a_i, r_i, s'_i) using the current policy
12:     for b = 1 ... B do
13:         Sample a batch, b, of experiences from the experience replay memory
14:         for u = 1 ... U do
15:             for i = 1 ... N do
16:                 # Calculate target Q-values for each example
17:                 y_i = r_i + δ_{s'_i} γ Q^{π_φ}(s'_i, max_{a'_i} Q^{π_θ}(s'_i, a'_i)) where δ_{s'_i} = 0
                        ↪   if s'_i is terminal, 1 otherwise
18:             end for
19:             # Calculate the loss, for example using MSE
20:             L(θ) = (1/N) Σ_i (y_i - Q^{π_θ}(s_i, a_i))^2
21:             # Update the network's parameters
22:             θ = θ - α ∇_θ L(θ)
23:         end for
24:     end for
25:     Decay τ
26:     if (m mod F) == 0 then
27:         # Update the target network
28:         φ = θ
29:     end if
30: end for
```

Line 10: **for** $m = 1 \ldots MAX\_STEPS$ **do**

Line 11: Gather and store $h$ experiences $(s_i, a_i, r_i, s'_i)$ using the current policy

Line 12: **for** $b = 1 \ldots B$ **do**

Line 13: Sample a batch, $b$, of experiences from the experience replay memory

Line 14: **for** $u = 1 \ldots U$ **do**

Line 15: **for** $i = 1 \ldots N$ **do**

Line 16: # Calculate target $Q$-values for each example

Line 17: $y_i = r_i + \delta_{s'_i} \gamma Q^{\pi_\varphi}(s'_i, \max_{a'_i} Q^{\pi_\theta}(s'_i, a'_i))$ where $\delta_{s'_i} = 0$
↪   if $s'_i$ is terminal, 1 otherwise

Line 18: **end for**

Line 19: # Calculate the loss, for example using MSE

Line 20: $L(\theta) = \frac{1}{N} \sum_i (y_i - Q^{\pi_\theta}(s_i, a_i))^2$

Line 21: # Update the network's parameters

Line 22: $\theta = \theta - \alpha \nabla_\theta L(\theta)$

Line 23: **end for**

Line 24: **end for**

Line 25: Decay $\tau$

Line 26: **if** $(m \mod F) == 0$ **then**

Line 27: # Update the target network

Line 28: $\varphi = \theta$

Line 29: **end if**

Line 30: **end for**