**Algorithm 1** Generic Q-learning (includes FQI and DQN as special cases

---

1: initialize $\phi_0$
2: initialize $\pi_0(\mathbf{a}|\mathbf{s}) = \epsilon\mathcal{U}(\mathbf{a}) + (1 - \epsilon)\delta(\mathbf{a} = \arg\max_{\mathbf{a}} Q_{\phi_0}(\mathbf{s}, \mathbf{a}))$ ▷ Use $\epsilon$-greedy exploration
3: initialize replay buffer $\mathcal{D} = \emptyset$ as a ring buffer of fixed size
4: initialize $\mathbf{s} \sim d_0(\mathbf{s})$
5: **for** iteration $k \in [0, \dots, K]$ **do**
6:     **for** step $s \in [0, \dots, S - 1]$ **do**
7:         $\mathbf{a} \sim \pi_k(\mathbf{a}|\mathbf{s})$                   ▷ sample action from exploration policy
8:         $\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})$                   ▷ sample next state from MDP
9:         $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}, \mathbf{a}, \mathbf{s}', r(\mathbf{s}, \mathbf{a}))\}$   ▷ append to buffer, purging old data if buffer too big
10:     **end for**
11:     $\phi_{k,0} \leftarrow \phi_k$
12:     **for** gradient step $g \in [0, \dots, G - 1]$ **do**
13:         sample batch $B \subset \mathcal{D}$                 ▷ $B = \{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_t)\}$
14:         estimate error $\mathcal{E}(B, \phi_{k,g}) = \sum_i \left(Q_{\phi_{k,g}} - (r_i + \gamma\max_{\mathbf{a}'} Q_{\phi_k}(\mathbf{s}', \mathbf{a}'))\right)^2$
15:         update parameters: $\phi_{k,g+1} \leftarrow \phi_{k,g} - \alpha\nabla_{\phi_{k,g}}\mathcal{E}(B, \phi_{k,g})$
16:     **end for**
17:     $\phi_{k+1} \leftarrow \phi_{k,G}$                     ▷ update parameters
18: **end for**

---