**Algorithm 1** Q-$\lambda$ learning

---

1: **function** QLAMBDALEARNING($\alpha, \gamma, \lambda$)
2:     Initialize $\theta$, hidden state $x_0$ and belief state $b_0$
3:     $t \leftarrow 0$
4:     $a_0, a_0^* \leftarrow$ CHOOSEACTION($\theta, b_0$)
5:     $\xi_0 \leftarrow$ GRIDWEIGHTS($b_0, a_0^*$)
6:     $e \leftarrow \mathbf{0}$
7:     **loop**
8:         Observe reward $r_t$ and observation $y_{t+1}$
9:         $b_{t+1} \leftarrow$ UPDATEBELIEFSTATE($b_t, a_t, y_{t+1}$)
10:         $a_t, a_t^* \leftarrow$ CHOOSEACTION($\theta, b_t$)
11:         **if** $a_t \neq a_t^*$ **then**
12:             $e \leftarrow \mathbf{0} + \xi_t$
13:         **else**
14:             $e \leftarrow e + \xi_t$
15:         $\xi_{t+1} \leftarrow$ GRIDWEIGHTS($b_t, a_t$)
16:         $\xi_{t+1}^* \leftarrow$ GRIDWEIGHTS($b_t, a_t^*$)
17:         $\delta \leftarrow r_t + \gamma \cdot \theta^T \xi_{t+1}^* - \theta^T \xi_t$
18:         $\theta \leftarrow \theta + \alpha \cdot \delta \cdot e$
19:         $e \leftarrow \gamma \cdot \lambda \cdot e$
20:         $t \leftarrow t + 1$

21: **function** CHOOSEACTION($\theta, b$)                          ▷ $\epsilon$-greedy
22:     $a^* \leftarrow \arg\max_a \theta^T[\text{GRIDWEIGHTS}(b, a)]$
23:     **if** Uniform(0,1) $\leq \epsilon$ **then**
24:         $a \leftarrow$ random action
25:     **else**
26:         $a \leftarrow a^*$
27:     **return** $a, a^*$

28: **function** GRIDWEIGHTS($b, a$)
29:     $\xi \leftarrow \mathbf{0}$
30:     **for** each particle $p \in b$ **do**
31:         $\xi_a \leftarrow \xi_a + \text{INTERPOLANTS}(p)$
32:     $\xi \leftarrow \frac{\xi}{\sum_{a,i}(\xi_{ai})}$
33:     $\xi_{a\ i+1} \leftarrow \text{var}(\xi)$
34:     **return** $\xi$

---

1

**Algorithm 2** Modified Upper Confidence Bound for Trees

---

1: **function** SELECTACTION($b, d$)
2:     Initialize $Q$ and $N$
3:     $s \leftarrow$ random state from belief state $b$
4:     **loop**
5:         SIMULATE($s, d$)
6:     **return** $\arg\max_a Q(a)$

7: **function** SIMULATE($s, d$)
8:     **if** $d = 0$ **then**
9:         **return** 0
10:    **if** $S \notin T$ **then**
11:        **for** $a \in A(s)$ **do**
12:            $(N(a), Q(a)) \leftarrow (N_0(a), Q_0(a))$
13:        $T = T \cup \{s\}$
14:        **return** ROLLOUT($s, d$)
15:    $a \leftarrow$ CHOOSEACTION($s$)
16:    $(s', r) \sim G(s, a)$
17:    $q \leftarrow r + \gamma$SIMULATE($s', d - 1$)
18:    $N(a) \leftarrow N(a) + 1$
19:    $Q(a) \leftarrow Q(a) + \frac{q - Q(a)}{N(a)}$

20: **function** ROLLOUT($s, d$)
21:    **if** $d = 0$ **then**
22:        **return** 0
23:    $a \leftarrow$ random action $a \in A$
24:    $(s', r) \sim G(s, a)$
25:    **return** $r + \gamma$ROLLOUT($s', d - 1$)

26: **function** CHOOSEACTION($s$)
27:    **for** $a \in A$ **do**
28:        **if** $N(s, a) = 0$ **then return** a
29:    **return** $\arg\max_a Q(a) + c\sqrt{\frac{\log(\sum_a N(a))}{N(a)}}$

---