

---

**Algorithm 1** HumanPolicyGenerator

---

**Input:** POLICYFUNCTION**Output:**  $\pi_{\text{POLICYFUNCTION}}^{(2)}$ 

- 1: Start from  $\pi_{\text{POLICYFUNCTION}}^{(2)}(s) = \text{NaN}, \forall s \in S$
  - 2: **for**  $s \in S$  **do**
  - 3:      $\pi_{\text{POLICYFUNCTION}}^{(2)}(s) = \text{POLICYFUNCTION}(s)$
  - 4: **end for**
- 

**Algorithm** POLICYFUNCTION:TABLE

---

**Input:**  $s \in S$ **Output:**  $a \in A$ 

---

**Algorithm** POLICYFUNCTION:BELOW

---

**Input:**  $s \in S, \{\text{LEFT}, \text{RIGHT}\}, \{\text{LEFT}, \text{RIGHT}\}$ **Output:**  $a \in A$ 

---

---

**Algorithm 2** ValueIteration

---

**Input:**  $\pi^{(2)}$ **Output:**  $\pi^{(1)}, V(\pi^{(1)}, \pi^{(2)})$ 

---

---

**Algorithm 3**  $Q$  learning with UCB exploration

---

**Input:**  $\varepsilon, \varepsilon_2, \delta, c_2$ , and `ITERS`.

**Output:**  $Q: S \times A \rightarrow [0, 1]$ .

- 1:  $Q_0(s, a), \hat{Q}_0(s, a) \leftarrow \frac{1}{1-\gamma}, N(s, a) = 0, R = \lceil \log(\frac{1}{\varepsilon_2(1-\gamma)}) / (1-\gamma) \rceil$ .
  - 2:  $L = \lfloor \log_2 R \rfloor, \varepsilon_L = \frac{1}{2^{L+2}} \varepsilon_2 (\log(1/(1-\gamma)))^{-1}$ .
  - 3:  $M = \max\{\lceil 2\log_2(\frac{1}{\varepsilon_L(1-\gamma)}) \rceil, 10\}, \varepsilon_1 \leftarrow \frac{\varepsilon}{24RM\log \frac{1}{1-\gamma}}, H = \frac{\log 1/((1-\gamma)\varepsilon_1)}{\log 1/\gamma}$ .
  - 4:  $\iota(k) \triangleq \log(SA(k+1)(k+2)/\delta), \alpha_k \triangleq \frac{H+1}{H+k}$ .
  - 5: **for**  $i = 1$ : `ITERS` **do**
  - 6:    $a_i \leftarrow \arg \max_{a'} \hat{Q}(s_i, a')$ .
  - 7:   Receive reward  $r_i$  and transit to  $s_{i+1}$ .
  - 8:    $N(s_i, a_i) \leftarrow N(s_i, a_i) + 1$ .
  - 9:    $k \leftarrow N(s_i, a_i), b_k = \frac{c_2}{1-\gamma} \sqrt{\frac{H\iota(k)}{k}}$
  - 10:    $Q_{i+1}(s, a) \leftarrow (1-\alpha)Q_i(s, a) + \alpha[r(s, a) + \gamma \max_{a'} \hat{Q}_{i+1}(s', a') + b_k]$ .
  - 11:    $\hat{Q}_{i+1}(s, a) \leftarrow \min[\hat{Q}_{i+1}(s, a), Q_i(s, a)]$ .
  - 12: **end for**
- 

---

**Algorithm 4** Maximize to Explore

---

**Input:**  $\eta$ , Hypothesis class  $\mathcal{H}$ , `VALUEITERATION`.

**Output:**  $\pi^{(1)}$

- 1:  $\mathcal{D} = \phi$
  - 2: **for**  $i = 1$ : `ITERS` **do**
  - 3:    $\hat{\pi} = \arg \max_{\pi^{(2)} \in \mathcal{H}} (V(\text{VALUEITERATION}(\pi^{(2)}), \pi^{(2)}) - \eta \mathcal{L}^{\pi^{(2)}}(\mathcal{D}))$
  - 4:    $\pi^{(1)} \leftarrow \text{VALUEITERATION}(\hat{\pi})$
  - 5:   Simulate with  $\pi^{(1)}$  until the game terminates and store  $(s, a, s')$  pairs into  $\mathcal{D}$ .
  - 6: **end for**
- 

---

**Algorithm 5** Upper Confidence Bound

---

**Input:**  $\alpha$ , Hypothesis class  $\mathcal{H}$ , `VALUEITERATION`.

**Output:**  $\pi^{(1)}$ .

- 1:  $N(\pi) \leftarrow 1$  and  $V(\pi) \leftarrow 0$  for all  $\pi$ .
  - 2: **for**  $i = 1$ : `ITERS` **do**
  - 3:    $\hat{\pi} = \arg \max_{\pi^{(1)} \in \text{VALUEITERATION}(\mathcal{H})} (V(\pi^{(1)}) + \sqrt{\frac{\alpha \log i}{N(\pi^{(1)})}})$ .
  - 4:    $\pi^{(1)} \leftarrow \hat{\pi}$
  - 5:    $N(\pi^{(1)}) \leftarrow N(\pi^{(1)}) + 1$ .
  - 6:   Simulate with  $\pi^{(1)}$  and receive the cumulative reward  $V$ .
  - 7:    $V(\pi^{(1)}) \leftarrow V$ .
  - 8: **end for**
-