
Algorithm 1 RUN-EPISODE (Q, π^h, π^r)

INPUT: State space S , Action space A

- 1: Initialize environment to start state; $R \leftarrow 0$; $i \leftarrow 0$;
 - 2: **loop**
 - 3: **if** goal state **then**
 - 4: **return** Q -functions, R
 - 5: **end if**
 - 6: Sample a_h from π^h and a_r from π^r
 - 7: Take action (a_h, a_r) and observe r, \mathbf{s}'
 - 8: Update Q -functions as in Eq. 2
 - 9: $R \leftarrow R + \gamma^i * r$
 - 10: $i \leftarrow i + 1$
 - 11: **end loop**
-