

# Reinforcement Learning Exercise

## Exercises 03

### 1 Dynamic programming

(a)

$$\begin{aligned} q_\pi(s, a) &= \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] \\ &= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\ &= \mathbb{E}_\pi[R_{t+1} + \gamma \sum_{s', a'} q_\pi(s', a') \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r \mid s, a) \left( r + \gamma \sum_{a'} \pi(a' \mid s') q_\pi(s', a') \right) \end{aligned}$$

$$\begin{aligned} q_{k+1}(s, a) &= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r \mid s, a) \left( r + \gamma \sum_{a'} \pi(a' \mid s') q_k(s', a') \right) \end{aligned}$$