

## Wprowadzenie do RL 5

Wszystkie zadania dotyczą środowiska **Frozen Lake**

### Zadanie 1

Napisz funkcję, która ze znajomości  $V(s)$  dla każdego stanu  $s$ , pozwala wyliczyć wartości zwrotów dla wszystkich akcji w pewnym stanie  $s$ , czyli  $Q(s, a)$ , dla każdej akcji  $a$  możliwej do wykonania w stanie  $s$ . Wykorzystaj formułę:

$$q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

Zapoznaj się z notatnikiem [FrozenLake\\_3.ipynb](#) (otwórz go w Colab). Znajdziesz tam wskazówki. Wykonaj podane tam polecenia.

### Zadanie 2

Napisz program implementujący poniższy **algorytm iteracji polityki** do środowiska *Frozen Lake* w celu znalezienia optymalnej polityki  $\pi_*$ . Przyjmij, że polityka  $\pi$  jest **deterministyczna**.

#### Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization  
 $V(s) \in \mathbb{R}$  and  $\pi(s) \in \mathcal{A}(s)$  arbitrarily for all  $s \in \mathcal{S}$
2. Policy Evaluation  
Loop:  
     $\Delta \leftarrow 0$   
    Loop for each  $s \in \mathcal{S}$ :  
         $v \leftarrow V(s)$   
         $V(s) \leftarrow \sum_{s', r} p(s', r | s, \pi(s)) [r + \gamma V(s')]$   
         $\Delta \leftarrow \max(\Delta, |v - V(s)|)$   
until  $\Delta < \theta$  (a small positive number determining the accuracy of estimation)
3. Policy Improvement  
     $policy\_stable \leftarrow true$   
    For each  $s \in \mathcal{S}$ :  
         $old\_action \leftarrow \pi(s)$   
         $\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$   
        If  $old\_action \neq \pi(s)$ , then  $policy\_stable \leftarrow false$   
    If  $policy\_stable$ , then stop and return  $V \approx v_*$  and  $\pi \approx \pi_*$ ; else go to 2

Zapoznaj się z notatnikiem [FrozenLake\\_4.ipynb](#) (otwórz go w Colab). Znajdziesz tam wskazówki. Wykonaj podane tam polecenia.