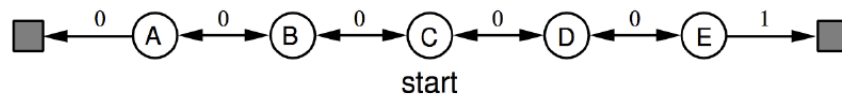


## Wprowadzenie do RL 8

### Zadanie 1

Dla spaceru przypadkowego omówionego na wykładzie:



Znajdź wartości  $V(A)$ ,  $V(B)$ ,  $V(C)$ ,  $V(D)$ ,  $V(E)$ :

- A. Wykorzystując programowanie dynamiczne (**równanie Bellmana**)
- B. Wykorzystując **algorytm TD(0)** (uzupełnij notatnik: [RandomWalk\\_TODO.ipynb](#))

Przyjmij, że wykorzystana jest **polityka stochastyczna** przy której prawdopodobieństwo ruchu w lewo i w prawo = 0.5.

### Zadanie 2

Napisz program implementujący poniższy **algorytm TD(0)** dla środowiska *FrozenLake* w celu wyliczenia polityki  $\pi$ .

#### Tabular TD(0) for estimating $v_\pi$

Input: the policy  $\pi$  to be evaluated

Algorithm parameter: step size  $\alpha \in (0, 1]$

Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

$A \leftarrow$  action given by  $\pi$  for  $S$

        Take action  $A$ , observe  $R, S'$

$V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

    until  $S$  is terminal

Wykorzystaj notatnik: [FrozenLake\\_TD\(0\)\\_TODO.ipynb](#)

### Zadanie 3

Napisz program implementujący poniższy **algorytm** dla środowiska *FrozenLake*.

#### Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$   
Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$   
Loop for each episode:  
  Initialize  $S$   
  Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)  
  Loop for each step of episode:  
    Take action  $A$ , observe  $R, S'$   
    Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)  
     $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$   
     $S \leftarrow S'; A \leftarrow A';$   
  until  $S$  is terminal

Wykorzystaj notatnik: [FrozenLake\\_SARSA\\_TODO.ipynb](#)

### Zadanie 4

Napisz program implementujący poniższy **algorytm** dla środowiska *FrozenLake*.

#### Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$   
Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$   
Loop for each episode:  
  Initialize  $S$   
  Loop for each step of episode:  
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)  
    Take action  $A$ , observe  $R, S'$   
     $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$   
     $S \leftarrow S'$   
  until  $S$  is terminal