

MDP

$\langle S, A, T, r, S_0, \gamma \rangle$

$$T(s, a, s') = \Pr[S_{t+1} = s' \mid S_t = s, A_t = a]$$

$$\sum_{s' \in S} T(s, a, s') = 1$$

Class MDPsim:

def estado_inicial():

ret S_0

def acciones_legales(s):

ret $A(s)$

def transición(s, a):

ret s'

def $r(s, a, s')$

ret recompensa

def terminal(s):

ret es_terminal

TD (0)

Estimar el valor de la política

```
def TD0 (mdp_sim,  $\pi$ ,  $\epsilon$ , num_episodios)
     $V = \{0 \text{ for } a \text{ in } \pi.keys()\}$ 
    for _ in range(num_episodios)
         $S = \text{mdp\_sim.estado\_inicial}()$ 
        for _ in range(num_iter)
            if mdp_sim.terminal(s)
                break
             $a = \pi(s)$ 
             $S' = \text{mdp\_sim.trans}(s)$ 
             $r = \text{mdp\_sim.reward}(s, a, s')$ 
             $V = V[s]$ 
             $V[s] = V[s] + \alpha \cdot (r + \text{mdp\_sim.}\gamma \times V[s'] - V[s])$ 
             $\Delta = \max(\Delta, \text{abs}(V - V[s]))$ 
             $S = S'$ 
        if  $\Delta < \epsilon$ 
            break
    return V
```