

# Sarsa Algorithm

$(s_t, a_t, r_t, s_{t+1}, a_{t+1})$

# Characteristics

Monte Carlo

Dynamic Programming

```
graph TD; MC[Monte Carlo] --> TD[Time Difference (TD)]; DP[Dynamic Programming] --> TD;
```

Time Difference (TD)

On Policy

# Algorithm

Initialize  $Q(s, a)$  arbitrarily

Repeat (for each episode):

Initialize  $s$

Choose  $a$  from  $s$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

Take action  $a$ , observe  $r, s'$

Choose  $a'$  from  $s'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -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

# Algorithm

Model

1	2	3	4	5	6
7	8	9	10	11	12
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30
31	32	33	34	35	36

Rewards

30 → 36: 2  
35 → 36: 2  
Else: 0

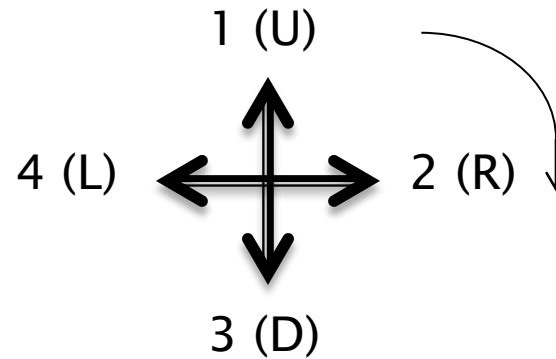
Why 0?

Because the algorithm does not know when and in which direction it has to go.

# Algorithm

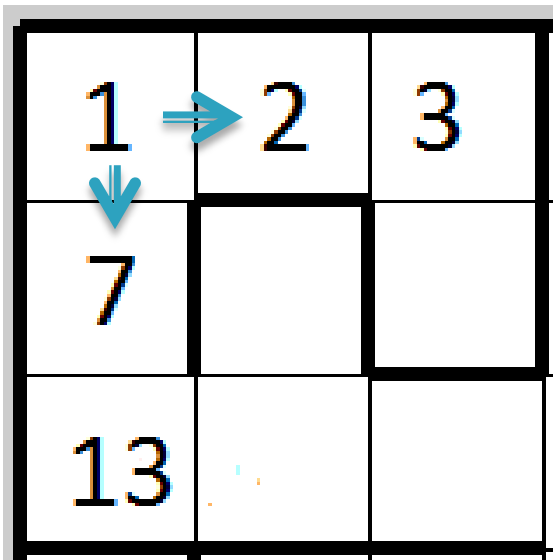
Structure  $Q(s,a)$

$X_1: U_1 \quad R_1 \quad D_1 \quad L_1$   
 $X_2: U_3 \quad R_2 \quad D_2 \quad L_2$   
 $X_3: U_2 \quad R_3 \quad D_3 \quad L_3$   
...  
...  
 $X_{36}: U_{36} \quad R_{36} \quad D_{36} \quad L_{36}$



# Algorithm

First iteration



S: state 1  
A: best action →  
R:  
S:  
B:

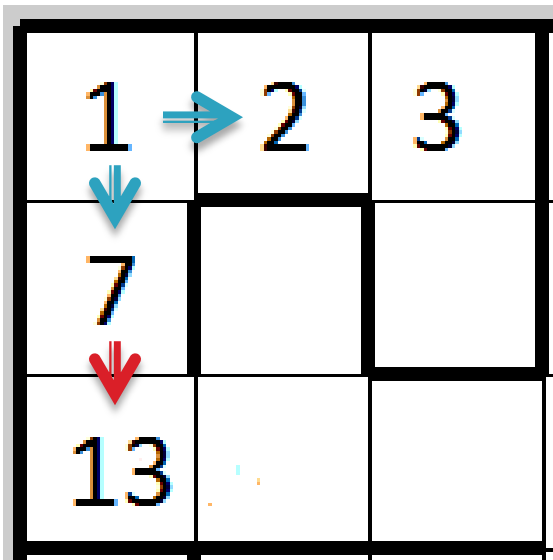
- $\epsilon$  prob. to choose a random action.
- if same reward, random action



Both have a reward of 0, so it will choose a random action.

# Algorithm

First iteration



S: state 1  
A: action 3 (down)  
R: reward 0  
S: state 7  
A: best action →

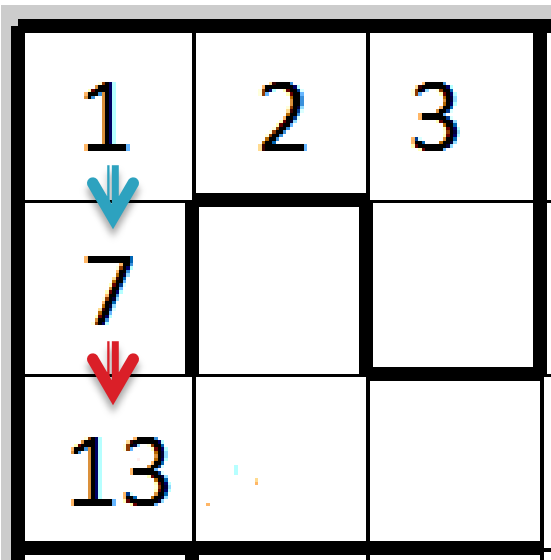
- $\epsilon$  prob. to choose a random action.
- if same reward, random action



Both have a reward of 0, so it will choose a random action.

# Algorithm

First iteration



S: state 1  
A: action 3 (down)  
R: reward 0  
S: state 7  
A: action 3 (down)



# Algorithm

First iteration

1	2	3
7		
13		

Algorithm

$$Q(s,a) \leftarrow Q(s,a) + \alpha[r + \gamma Q(s',a') - Q(s,a)]$$

$$Q(1,3) \leftarrow Q(1,3) + \alpha[0 + \gamma Q(7,3) - Q(1,3)]$$

$\alpha$ : learning rate

$r$ : reward

$\gamma$ : discount factor

# Algorithm

After many iterations

1	2	3	4	5	6
7	8	9	10	11	12
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30
31	32	33	34	35	36

Final  $Q(s,a)$

1:	0	0.0000	<b>0.0007</b>	0
7:	0.0002	0	<b>0.0033</b>	0
13:	0.0006	<b>0.0143</b>	0	0
14:	0.0007	<b>0.0276</b>	0	0.0001
15:	0	<b>0.0758</b>	0	0.0055
16:	0.0043	0	<b>0.3093</b>	0.0239
22:	0.0814	<b>0.4906</b>	0	0.0546
23:	0	0.0063	<b>0.9059</b>	0.1223
29:	0.2403	0	<b>1.3936</b>	0
35:	0.4974	<b>1.9444</b>	0	0

# Explanation

## Complementary explanation

<http://laid.delanover.com/reinforcement-learning-sarsa-algorithm-a-practical-case>

## Video

<https://www.youtube.com/watch?v=MZSK-Ho2-NA>

