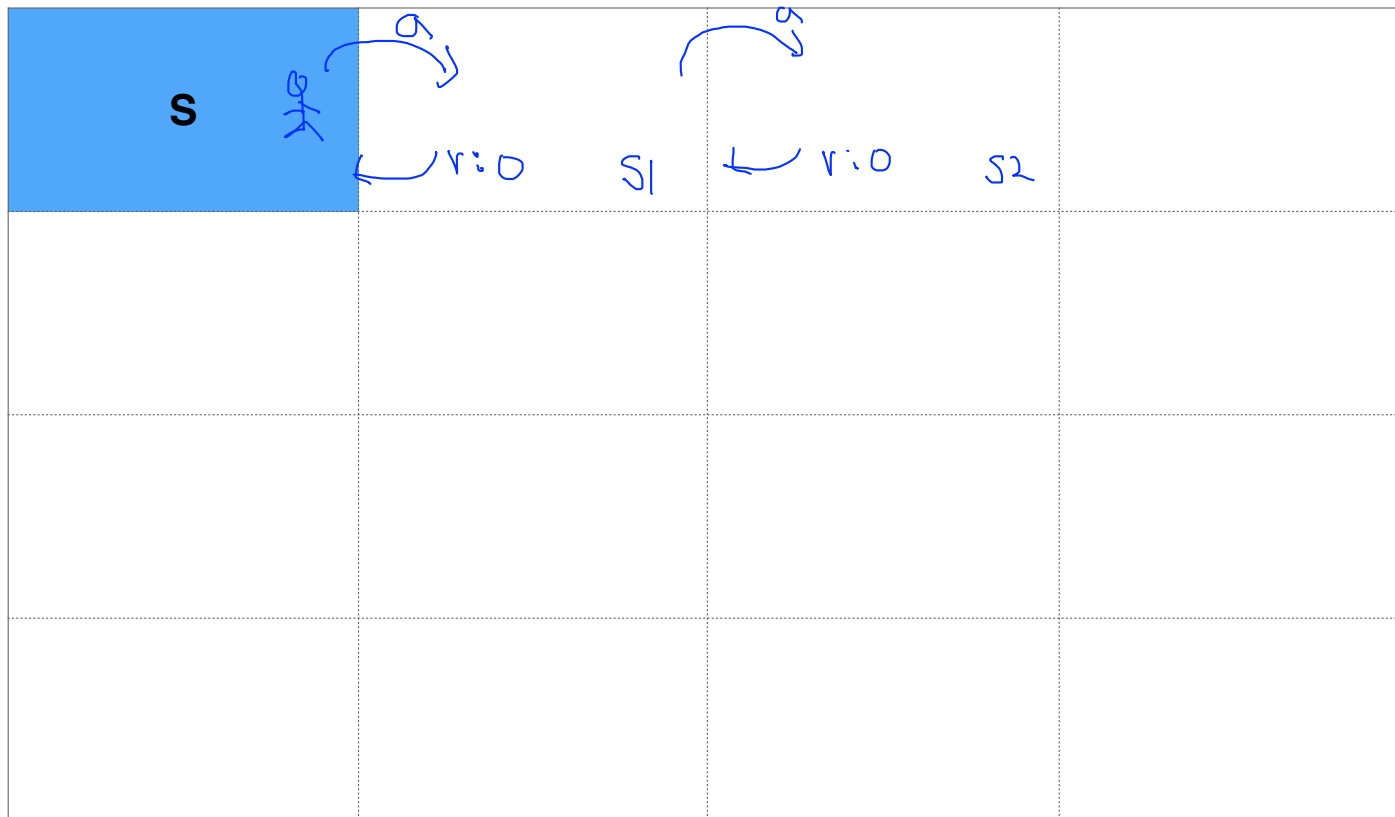




Lecture 3: Q-learning (table)

Reinforcement Learning with TensorFlow&OpenAI Gym
Sung Kim <hunkim+ml@gmail.com>

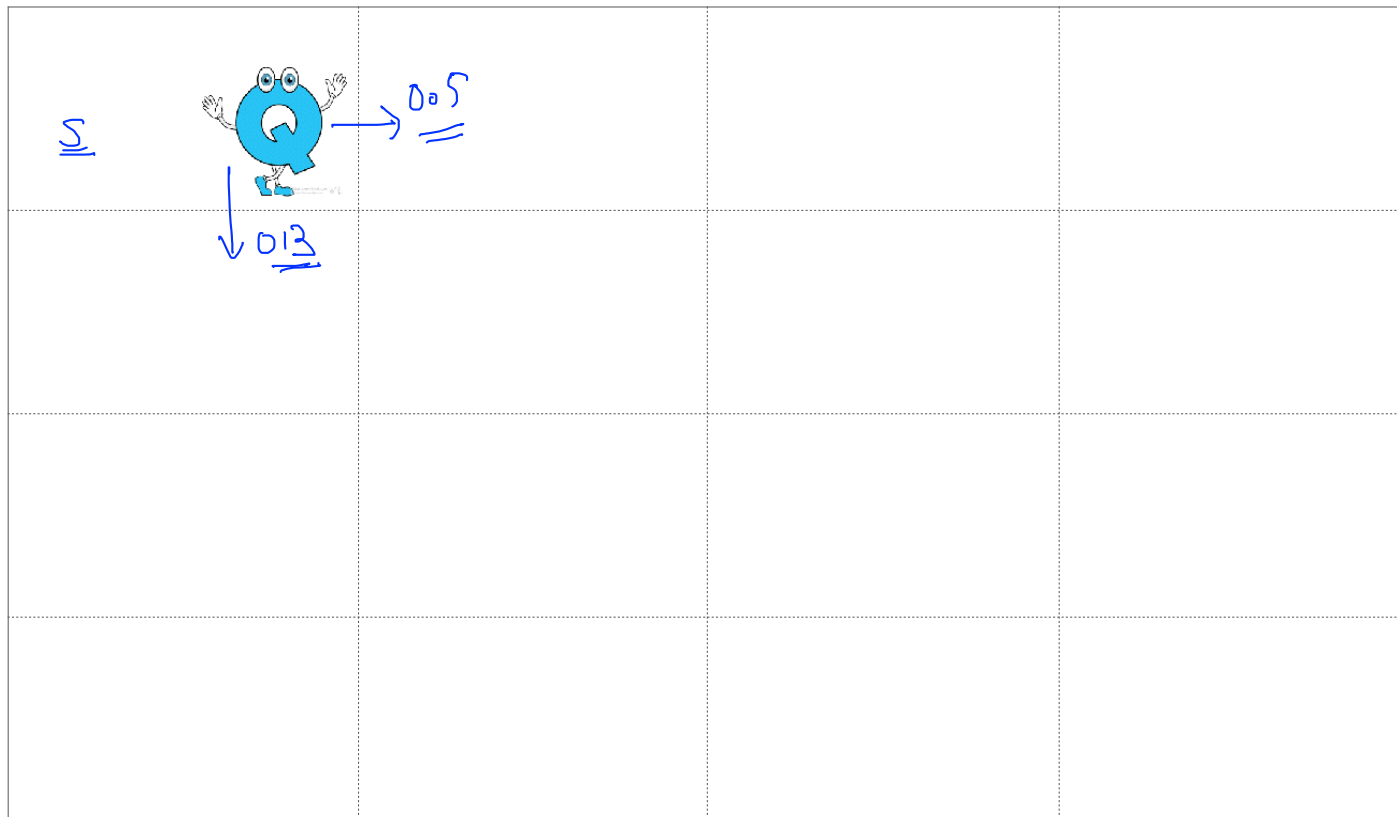
Try Frozen Lake, Real Game?



Frozen Lake: Random?

s			

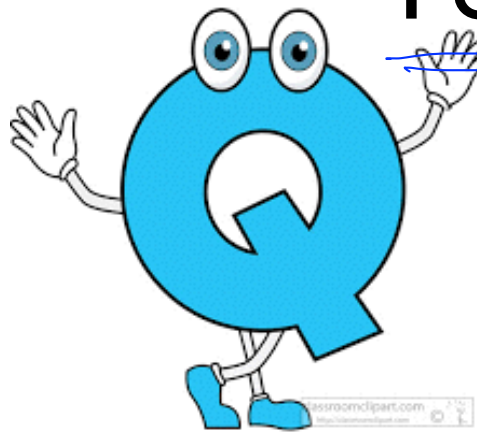
Frozen Lake: Even if you know the way, ask.
“아는 길도, 물어가라”



Q-function (state-action value function)



Q (state, action)



Policy using Q-function

Q (state, action)

Q (s1, LEFT): 0

Q (s1, RIGHT): 0.5

Q (s1, UP): 0

Q (s1, DOWN): 0.3

① max

$$\boxed{\max_a} Q(s1, a)$$

② max \rightarrow arg

RIGHT $\leftarrow \boxed{\arg \max_a} Q(s1, \underline{a})$

Optimal Policy, π and Max Q



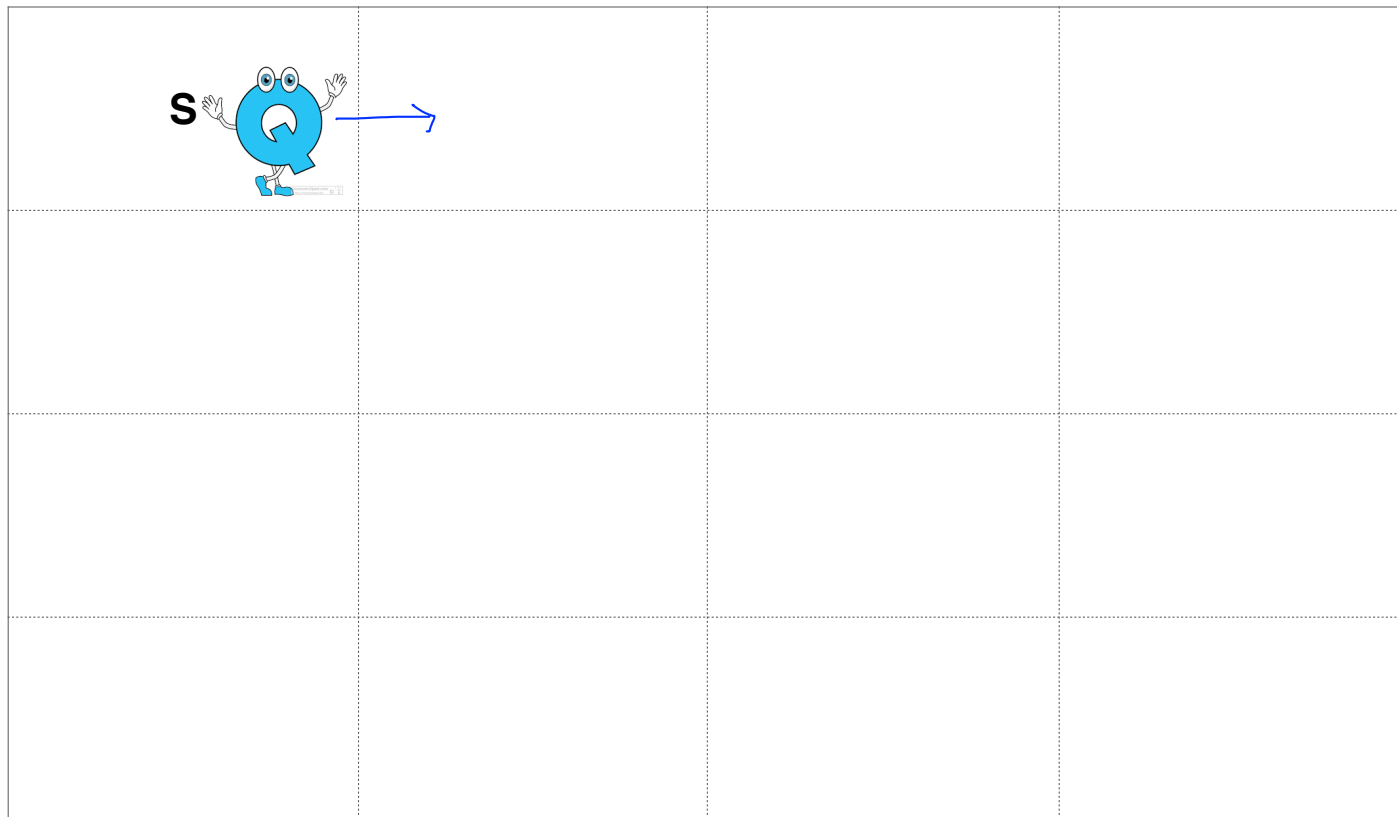
Q (state, action)



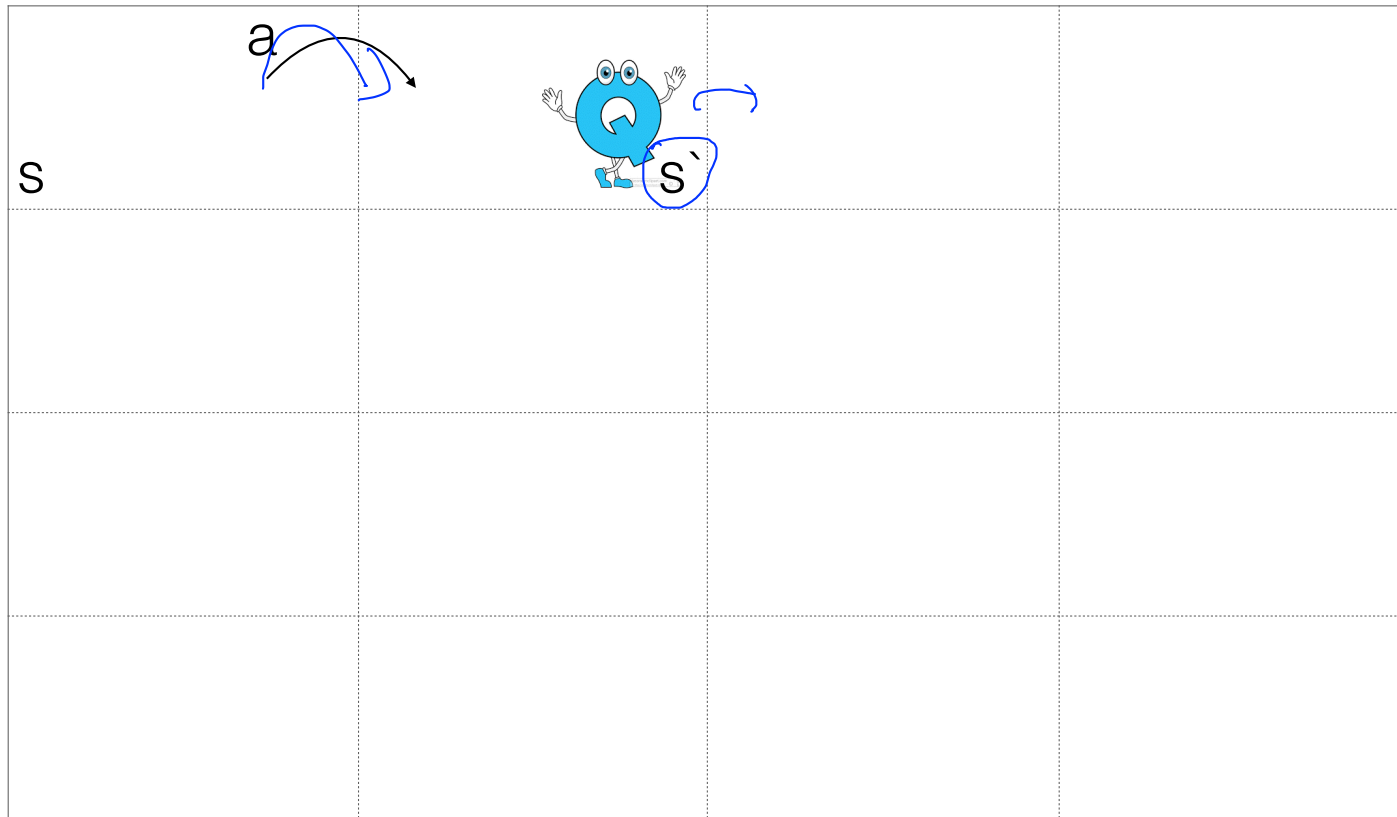
$$\text{Max } Q = \max_{a'} Q(s, a')$$

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} Q(s, a)$$

Frozen Lake: optimal policy with Q



Frozen Lake: optimal policy with Q



Finding, Learning Q

- Assume (believe) Q in s' exists!

- My condition

- I am in s
- when I do action a, I'll go to s'
- when I do action a, I'll get reward r
- Q in s', Q(s', a') exist!

$Q(s, a)$

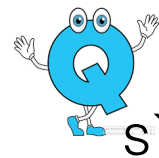
- How can we express $Q(s, a)$ using $Q(s', a')$?

$$\rightarrow = r + \max_{a'} Q(s', a')$$

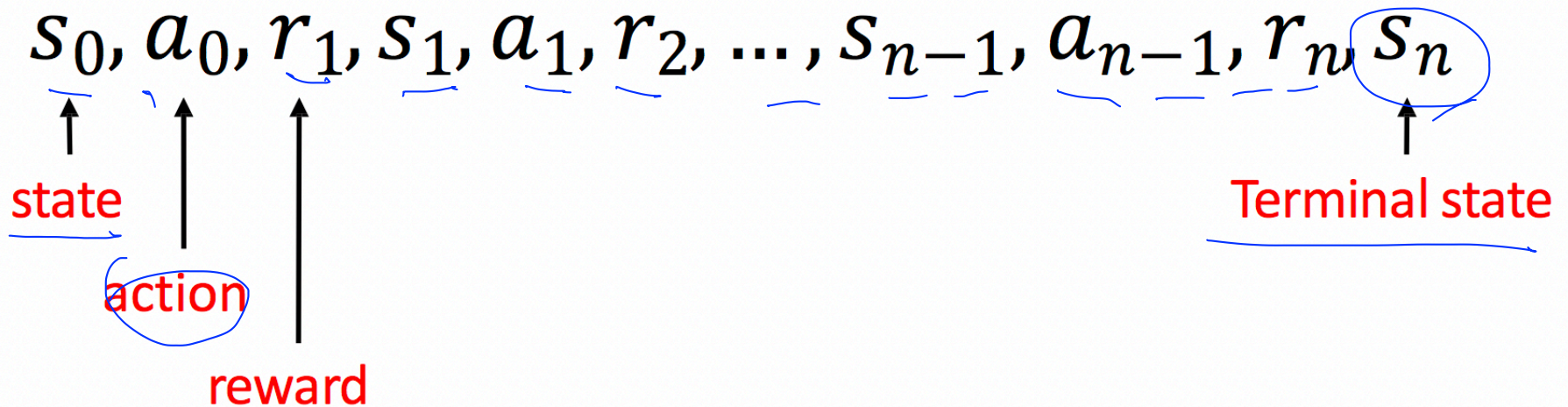
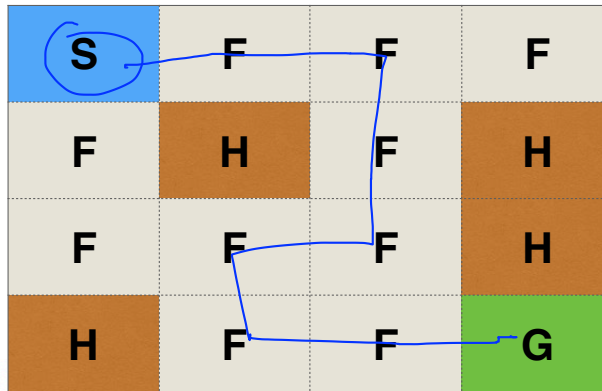


Learning Q (s, a)?

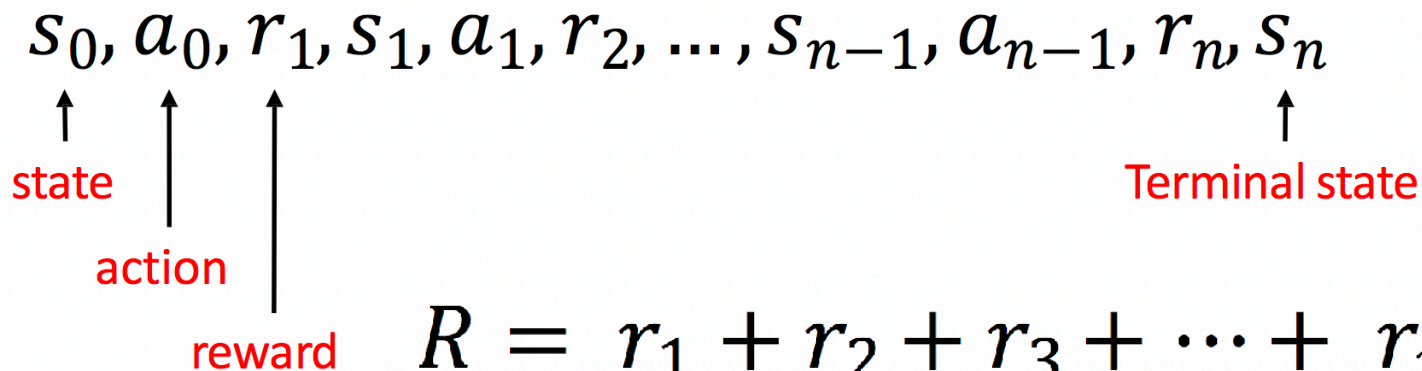
$$\begin{aligned} & \textcircled{S} \\ & Q(s, a)? \\ & = r + \max_{a'} Q(s', a') \end{aligned}$$



State, action, reward



Future reward



$$\underline{R} = \underline{r_1} + \underline{r_2} + \underline{r_3} + \dots + \underline{r_n} \quad * \text{optimal}$$

$$R(t) = r_t + R(t+1)$$

$$R_t = r_t + \boxed{r_{t+1} + r_{t+2} + \dots + r_n}$$

$$R_t^* = r_t + \max R(t+1)$$

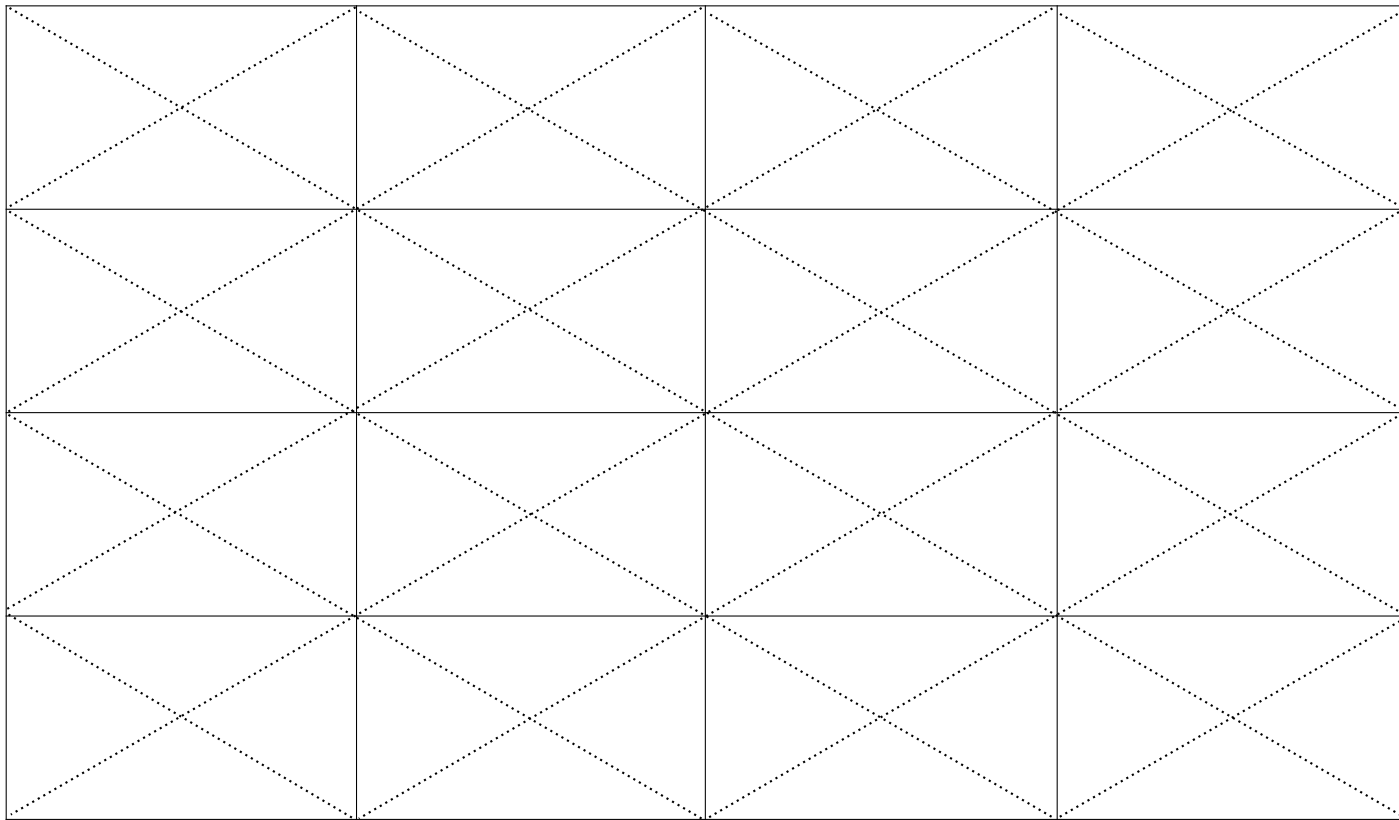
$$\underline{R(t+1)} = \underline{r_{t+1} + r_{t+2} + \dots + r_n} //$$

$$\boxed{Q(s, a) = r + \max_{a'} Q(s', a')}$$

Learning Q (s, a)?

$$\hat{Q}(s, a) \leftarrow r + \max_{a'} \hat{Q}(s', a')$$

Learning $Q(s, a)$: 16x4 Table
16 states and 4 actions (up, down, left, right)



Learning $Q(s, a)$: Table

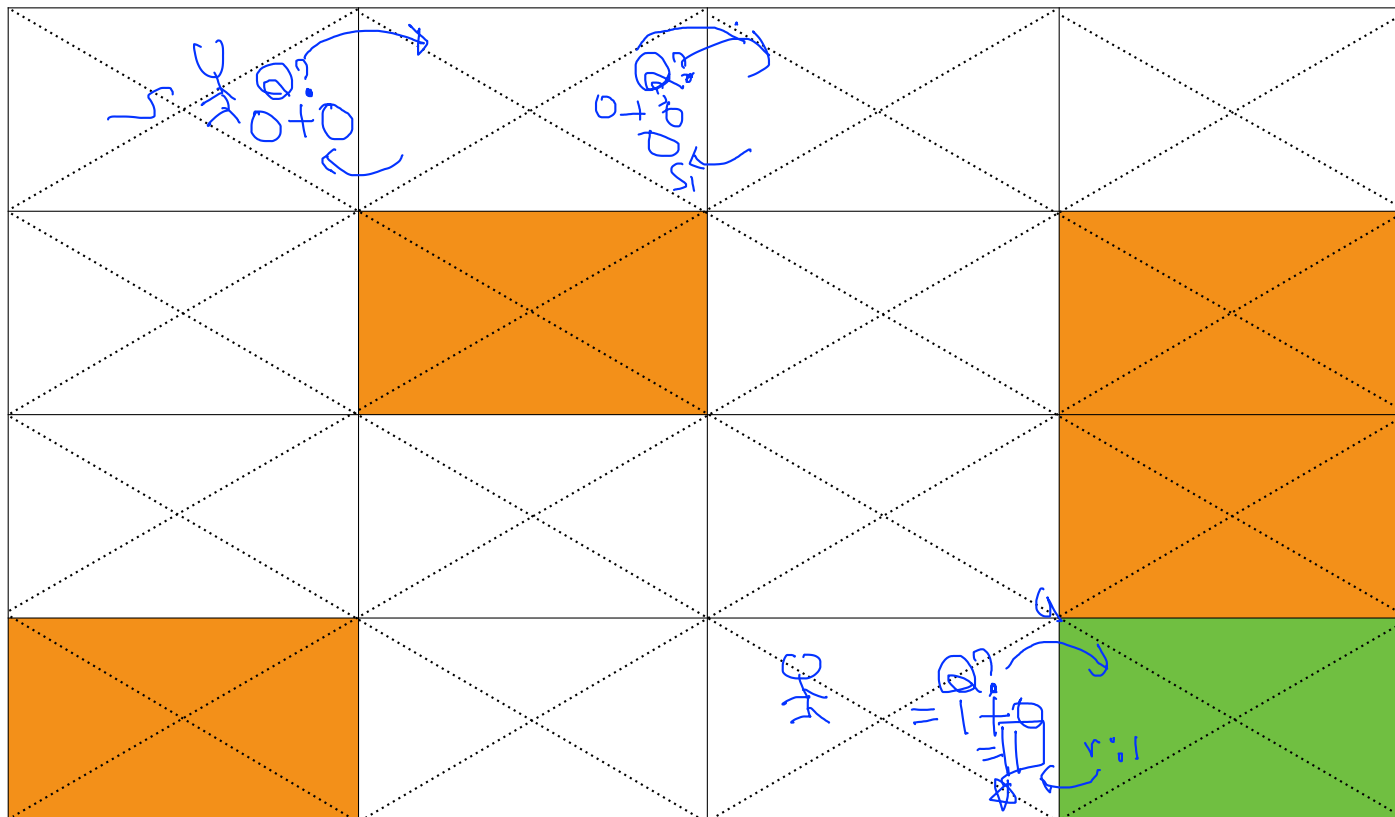
initial Q values are 0

A 4x4 grid of triangles. Each triangle contains the number '0'. The triangles are arranged in a regular pattern, with solid lines forming the outer grid and dotted lines forming the internal triangle boundaries.

$$= \left[r + \max_{a'} Q(s', a') \right]$$

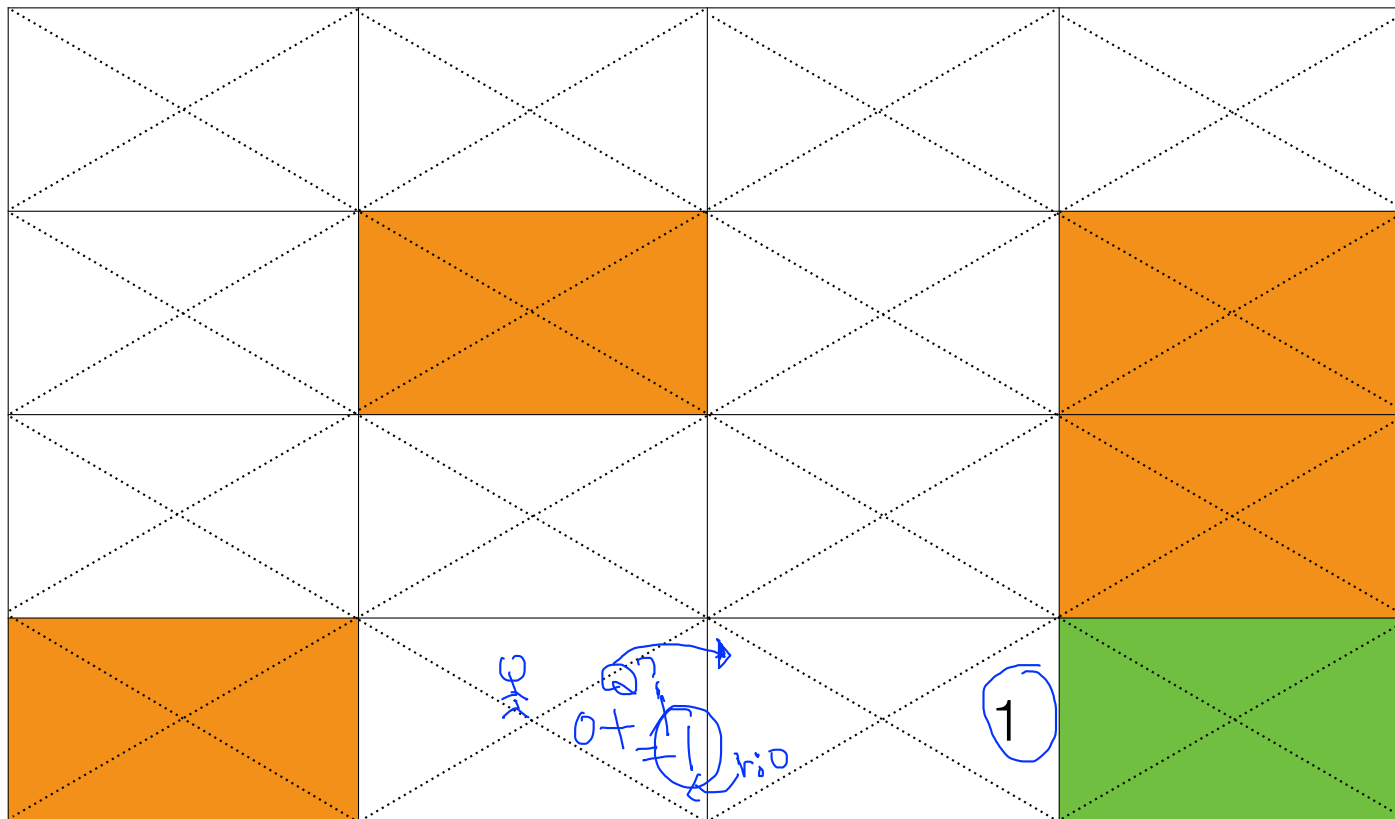
Learning $Q(s, a)$ Table (with many trials)

initial Q values are 0



Learning $Q(s, a)$ Table (with many trials)

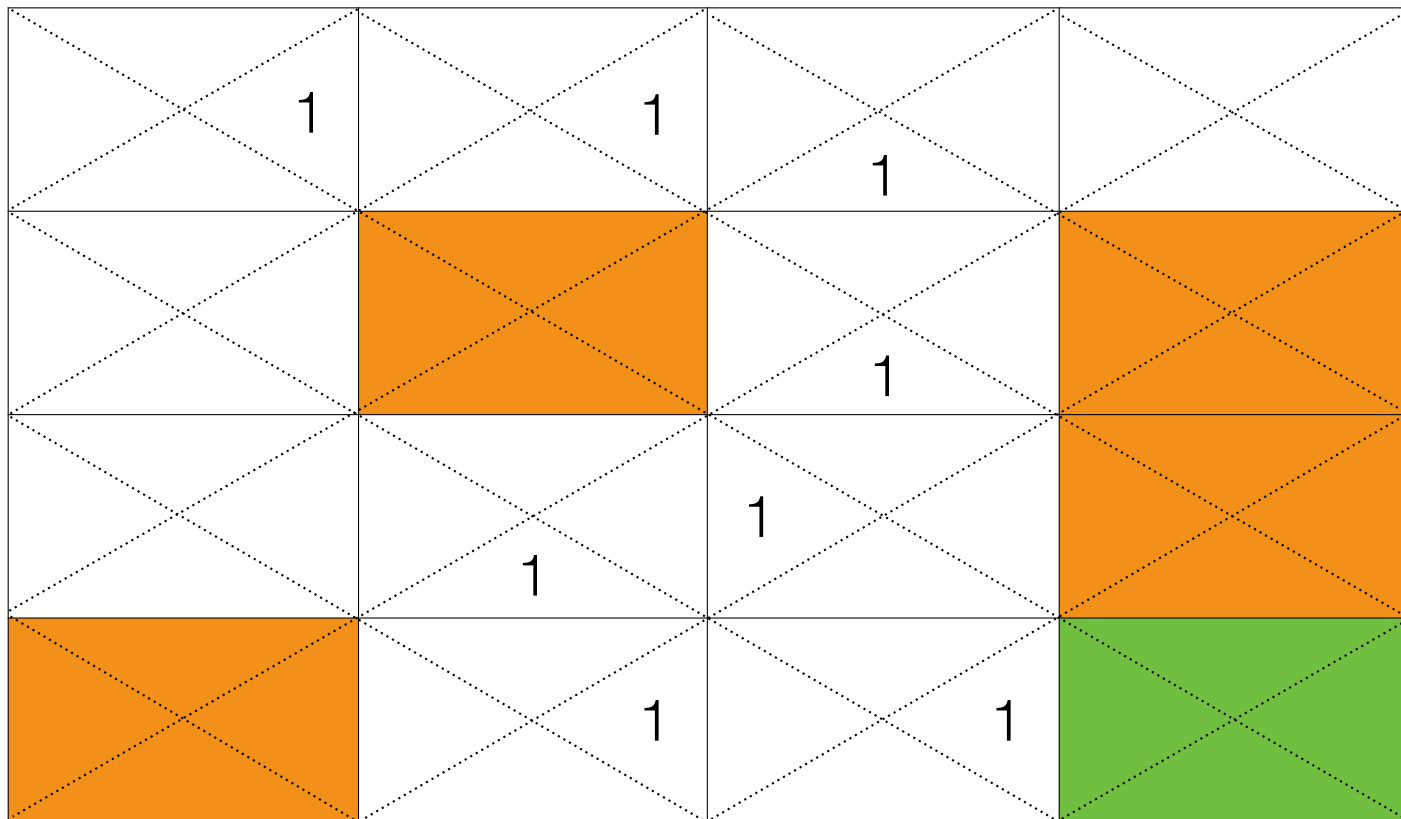
initial Q values are 0



$$Q(s_{14}, a_{\text{right}}) = r = 1$$

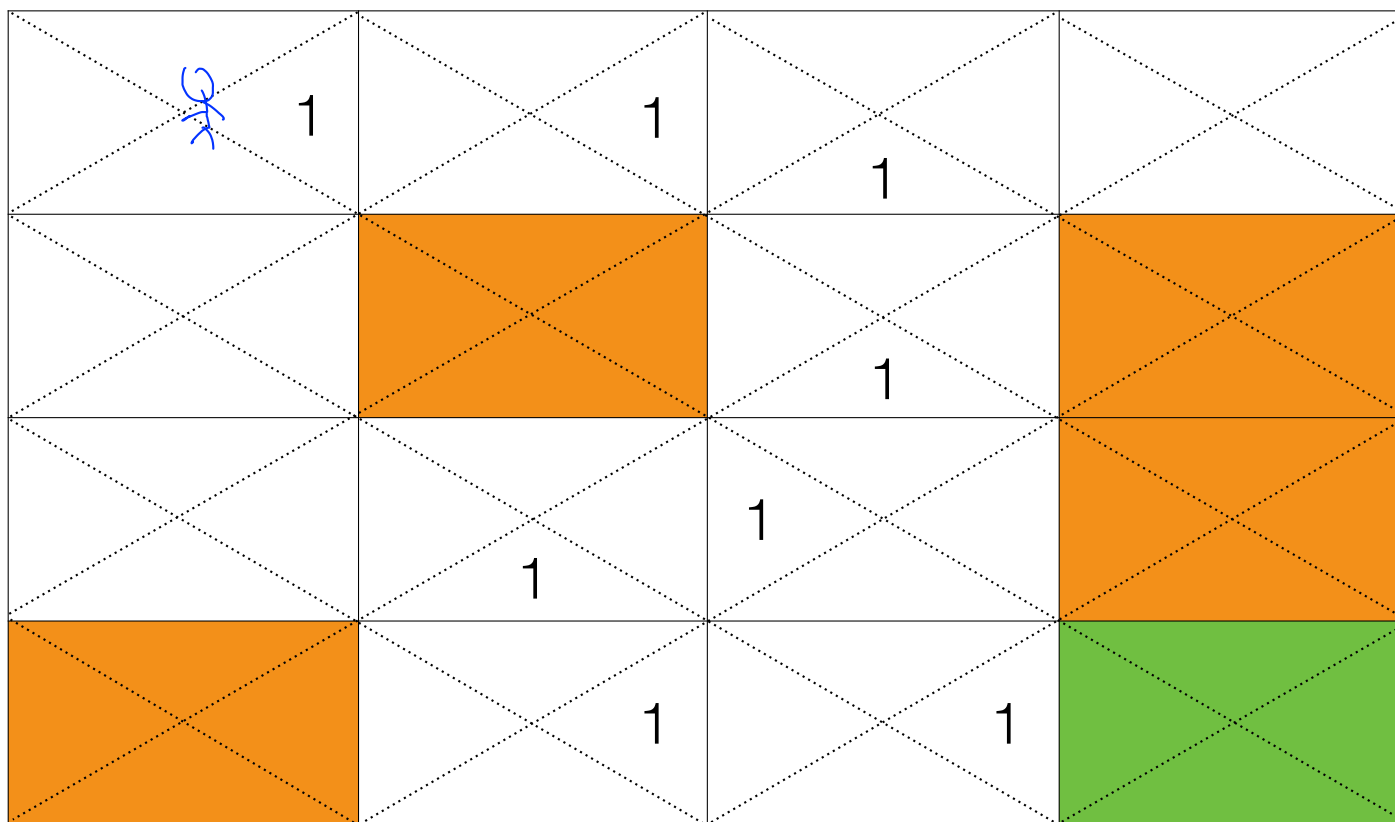
$$Q(\underline{s}_{13}, a_{\text{right}}) = r + \max(Q(s_{14}, a)) = 0 + \max(0, 0, 1, 0) = 1$$

initial Q values are 0



Learning $Q(s, a)$ Table: optimal policy

$$\pi^*(s) = \operatorname{argmax}_a Q(s, a)$$



Dummy Q-learning algorithm

For each s, a initialize table entry $\hat{Q}(s, a) \leftarrow \underline{\underline{0}}$

Observe current state s

Do forever:

- Select an action a and execute it
- Receive immediate reward r
- Observe the new state s'
- Update the table entry for $\hat{Q}(s, a)$ as follows:

$$\hat{Q}(\underline{s}, \underline{a}) \leftarrow \underline{r} + \max_{a'} \hat{Q}(s', a')$$

- $s \leftarrow s'$

Next
Lab: Dummy Q-learning Table

