# Conceptual View of Q-Learning

- Q-Learning is a form of Reinforcement Learning.

- It uses a random exploration paired with delayed reward to generate values for a "Q table".

- For some examples - the Q table is used to select the best action to move to a goal state given the agent is in a particular state

$$S(i) + a(j) \longrightarrow S(k)$$

$\uparrow$

state $(i)$ + Action $(j)$ $\rightarrow$ State $(k)$

results in

usually - in Q learning there are a finite number of states an agent can be in

there are some # of actions available to the agent. Basic systems would have all actions for each state

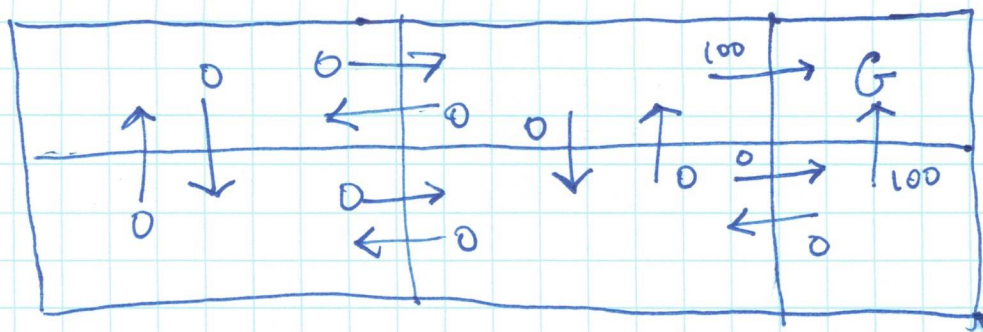new state as a result of action $(j)$

(1)



table shows immediate reward

- Only moving to the goal state provides a reward — that is an immediate reward.
- Every where else no immediate reward is given.

> A key part of Q learning is exploration of the state space.

- Agent is placed at a random state, and a random action is taken

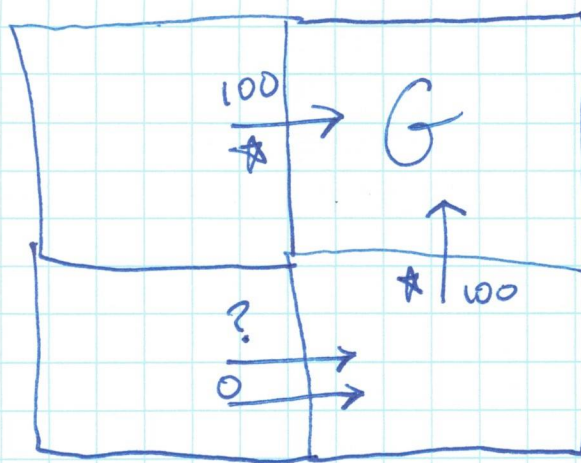  > it is the exploration that fills out the state table.

  filling out the table relies on the concept of delayed reward

  $$\underset{\text{reward}}{\text{delayed}} = \underset{\substack{\text{at state} \\ \uparrow \\ \text{current}}}{\text{imediate reward}} + \underset{\substack{\text{delay} \\ \text{factor}}}{\eta} * \underset{\substack{\text{immediate} \\ \text{successor state}}}{\text{reward of}}$$

G - goal
state

✳ — immediate reward for moving to the goal
state

? — here we need to assign reward based on
delay

? delay

O immediate

let $n = .9$

$$\text{delayed reward} = \text{immediate reward} + n * \text{reward of immediate successor state}$$

$$90 = 0 + .9 * 100$$

most reward available
from this state
( you have to pick
the correct action to
get it )

③

Basic Q Algorithm

for ( j=0 ; j < nEpochs ; j++ ) {

   myState = select Random State

  ~~●●●~~

   while ( myState != goalState ) {

      myAction = select Random Action

      ~~●●●●●●●●●●●●●●●●●~~

      nextState = ~~●●●●~~
        qMatrix adj [myState][myAction]

      immediate reward = immediate State [
      nextState].
        immediate
        reward

      delayed reward = immediate reward +
        $n$ * max ( all rewards
          from
          successor states )

q Table →

      stateTable [myState][myAction]

      ~~delayedReward =~~
        delayed Reward

      myState = nexState

  }   }

④

| 3 | 4 | 5 |
|---|---|---|
| 0 | 1 | 2 |

qImmyAction [ nStates ] =

[ 0, 0, 0, 0, 0, 0 ]

left, right, up, down

qAdj = [ [-1, 1, 3, -1],

qReward = [ [-1, 0, 0, -1],

qAdj [nStates]
[nActions]

[ 0, 2, 4, -1],

[ 1, -1, 5, -1],

[ -1, 4, -1, 0],

nStates = 6

[ 3, 5, -1, 1]

nActions = 4

[ 4, -1, -1, 2]]

qReward
[nStates]
[nActions]

[ 0, 0, 0, -1],

[ 0, -1, 0, -1],

[ -1, 0, -1, 0],

[ 0, 0, -1, 0]

[ 0, -1, -1, 0]]

goal = 5

lastState = nStates - 1

epochs = 0

qImmediate [goal] = 100

lastAction = nActions - 1

exploring = true

$R = .9$

while (exploring) {

   myState = random (0, lastState)

   while (myState != goal) {
      /* Get a valid random action. */
      myAction = random (0, lastAction)

      while ( qAdj [myState] [myAction] != -1) {

         myAction = random (0, lastAction)

     }

      nextState = qAdj [myState] [myAction]

```
            immediateReward = qImmediate [nextState]
            delayedReward = max ( qReward [nextState]) * ʀ
            totalReward = immediateReward + delayedReward
            qReward [myState][myAction] = totalReward
            myState = nextState
        }

    epochs ++

    if (epochs == maxEpochs) {
        exploring = false
    }
}
```
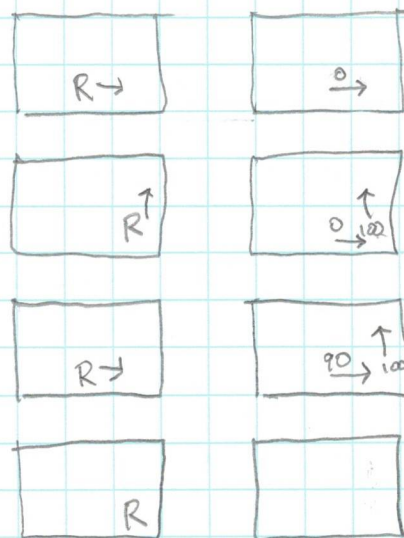
# Basic Q Algorithm

qImmediate     is annay that holds immeddiate reward given when agent moves to that state from an adjacent state

qMatrix     is an annay that holds the adjacent states given an action

qTable     holds the cummulitive rewards of recieved when moving upon one state to another

```
for (j=0; j < nEpochs; j++) {
    myState = selectRandomState();

    While (myState != goalState) {
        myAction = selectRandomAction();
        nextState = qMatrix[myState][myAction]
        immediateReward = qImmediate[next
                                      State]

        delayedReward = max(qTable
                                 [nextState]
                    * decay
```

```
total Reward = immediate Reward +
                  delayed Reward
qTable [ myState ][ myAction ] =
                      total Reward

myState = next state
} /* end while */

} /* end for */
```