# The link between an optimal Value and policy

In an optimal Q-Table we also have an optimal policy because we know what is the best action to take in each state.

$$\pi^*(s) = \arg\max_a Q^*(s, a)$$

But, in the beggining, our Q-table is useless since it gives arbitrary values.

As the agent explores the environment and we update the Q-table, it will give us better and better approximations to the optimal policy

## Q-Learning algorithm
Input: policy $\pi$, num-episodes, $\alpha$, GLIE $\{\epsilon_i\}$
Output: value function $Q$

Initialize Q-table (all ∅'s)

for $i \leftarrow 1$ to num-episodes do
  $\epsilon \leftarrow \epsilon_i$
  Observe $S_0$
  $t \leftarrow \emptyset$

  repeat
    Choose action $A_t$ using policy derived from $Q$ ($\epsilon$-greedy)
    Take action $A_t$ and observe $R_{t+1}$, $S_{t+1}$
    $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(...)$ (update state-action pair)
    $t \leftarrow t+1$
  until $S_t$ is terminal
  end
  return $Q$