

Algorithm 10: First-Visit GLIE MC Control

```
Input: positive integer num_episodes, GLIE \{\epsilon_i\}
Output: policy \pi (\approx \pi_* if num\_episodes is large enough)
Initialize Q(s, a) = 0 for all s \in S and a \in A(s)
Initialize N(s, a) = 0 for all s \in S, a \in A(s)
for i \leftarrow 1 to num_episodes do
   \epsilon \leftarrow \epsilon_i
    \pi \leftarrow \epsilon-greedy(Q)
    Generate an episode S_0, A_0, R_1, \dots, S_T using \pi
    for t \leftarrow 0 to T-1 do
        if (S_t, A_t) is a first visit (with return G_t) then
            N(S_t, A_t) \leftarrow N(S_t, A_t) + 1
            Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)}(G_t - Q(S_t, A_t))
    end
end
return \pi
```

Algorithm 11: First-Visit Constant- α (GLIE) MC Control

```
Input: positive integer num_episodes, small positive fraction \alpha, GLIE \{\epsilon_i\}
Output: policy \pi (\approx \pi_* if num\_episodes is large enough)
Initialize Q arbitrarily (e.g., Q(s, a) = 0 for all s \in S and a \in A(s))
for i \leftarrow 1 to num_episodes do
   \epsilon \leftarrow \epsilon_i
    \pi \leftarrow \epsilon-greedy(Q)
    Generate an episode S_0, A_0, R_1, ..., S_T using \pi
    for t \leftarrow 0 to T-1 do
        if (S_t, A_t) is a first visit (with return G_t) then
        Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(G_t - Q(S_t, A_t))
    end
end
return \pi
```

$$Q \leftarrow Q + N G - Q$$

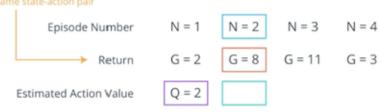
After each episode, we can calculate a new action value estimate

from the old action value estimate,

the most recently sampled return,

and the total number of first visits to the state-action pair.

Corresponding to the same state-action pair



$$Q \leftarrow Q + \frac{1}{N}G - Q$$