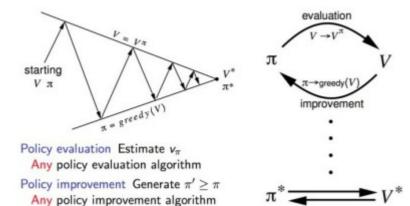
···



1. Initialization

 $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Repeat

$$\Delta \leftarrow 0$$

For each $s \in S$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta,|v - V(s)|)$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \theta$ (a small positive number)

3. Policy Improvement

policy- $stable \leftarrow true$

For each $s \in S$:

$$a \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \arg\max_{a} \sum_{s',r} p(s',r|s,a) \big[r + \gamma V(s') \big]$$

If $a \neq \pi(s)$, then policy-stable \leftarrow false

If policy-stable, then stop and return V and π ; else go to 2

```
Initialize array V arbitrarily (e.g., V(s) = 0 for all s \in S^+)
Repeat
    \Delta \leftarrow 0
    For each s \in S:
        v \leftarrow V(s)
         V(s) \leftarrow \max_{a} \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]
         \Delta \leftarrow \max(\Delta, |v - V(s)|)
until \Delta < \theta (a small positive number)
Output a deterministic policy, \pi, such that
    \pi(s) = \mathop{\arg\max}_{a} \sum_{s',r} p(s',r|s,a) \big[ r + \gamma V(s') \big]
```

```
初始化Q(s,a), \forall s \in S, a \in A(s),任意的数值,并且Q(terminal-state,\cdot)=0 重复(对每一节episode): 初始化 状态S 重复(对episode中的每一步): 使用某一个policy比如(\epsilon-greedy)根据状态S选取一个动作执行 执行完动作后,观察reward和新的状态S' Q(S_t,A_t) \leftarrow Q(S_t,A_t) + \alpha(R_{t+1} + \lambda \max_a Q(S_{t+1},a) - Q(S_t,A_t)) S \leftarrow S' 循环直到S终止
```