# Policy Improvement:

| S | A | $Q^{\hat{\pi}}(S,A) = R(S,A) + \gamma \sum_{S'} T(S,A,S') V^{\hat{\pi}}(S')$ |
|---|---|---|
| (0,0) | Right | $R([0,0], \text{Right}) + \gamma \sum [T([0,0], \text{Right}, [0,1]) V([0,1])$ |

$$+ T([0,0], \text{Right}, [1,0]) V([1,0])]$$

$$= 0 + 0.9 \times [0.9 \times 3.1 + 0.1 \times 0] = \boxed{2.511}$$

**Left** $\quad R([0,0], \text{Left}) + \gamma [T([0,0], \text{Left}, [0,0]) V([0,0])]$

$$= 0 + 0.9 \times [\,1 \times (-4.05)]$$

$$= -3.645$$

**Up** $\quad R([0,0], \text{Up}) + \gamma [T([0,0], \text{Up}, [0,0]) V([0,0])]$

$$= 0 + 0.9 \times [1 \times (-4.05)]$$

$$= -3.645$$

**Down** $\quad R([0,0], \text{Down}) + \gamma [T([0,0], \text{Down}, [1,0]) V([1,0])$

$$+ T([0,0], \text{Down}, [0,1]) V([0,1])]$$

$$= 0 + 0.9 \times [0.9 \times 0 + 0.1 \times 3.1]$$

$$= 0.279$$

**Nothing** $\quad R([0,0], \text{Nothing}) + \gamma [T([0,0], \text{Nothing}, [0,0]) V([0,0])]$

$$= 0 + 0.9 \times [1 \times (-4.05)]$$

$$= -3.645$$

In the new policy, $\pi([0,0]) = \text{Right}$ because the computed Q-value of action Right in state $[0,0]$ is the highest among all actions.

Do the same for all remaining states.