

$$Q(s,a) = Q(s,a) + \alpha [r + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$

Initial

	a1	a2
s1	0	0
s2	0	0
s3	0	0

$$\alpha = 0.5 \quad \gamma = 0.5$$

s: s1

a: a1

r: 0.5

s': s2

$$Q(s1, a1) = 0 + 0.5 [0.5 + 0.5(0) - 0] = 0.25$$

	a1	a2
s1	0.25	0
s2	0	0
s3	0	0

s: s2

a: a2

r: -1

s': s3

$$Q(s2, a2) = 0 + 0.5 [-1 + 0.5(0) - 0] = -0.5$$

	a1	a2
s1	0.25	0
s2	0	-0.5
s3	0	0

Optimal Policy

- From state 1, best action is a1 \rightarrow s2 because it yields a reward of 0.5 or 0 from s2 \rightarrow s1 by a1
- From s2, best action is a1 \rightarrow s1 because it has a reward of 0 ~~to~~ from s1 \rightarrow s2 through a loop. So you get rewards of 0 and 0.5 infinitely.
- From s3, best action is a2, taking you to s2. Yields a reward of 0. It also has a loop from s2 \rightarrow s1 resulting in a self loop which yields a constant reward of 0.5 and 0