

When it got reward r and the situation changed from s to s' after it took an action a , it revises the Q value for s and a as

$$\Delta Q(s, a) = \alpha \cdot \left(r + \gamma \max_{b \in A(s')} Q(s', b) - Q(s, a) \right)$$

under the learning rate α ($0 < \alpha < 1$) and discount rate γ ($0 < \gamma < 1$).