

An **off-policy learner** learns the value of the optimal policy independently of the agent's actions.

An **on-policy learner** learns the value of the policy being carried out by the agent, including the exploration steps.

SARSA: $Q(S_t, A_t) = Q(S_t, A_t) + \alpha [R(t+1) + \text{discount} * Q(S_{t+1}, \underline{A_{t+1}}) - Q(S_t, A_t)]$

Q-learn: $Q(S_t, A_t) = Q(S_t, A_t) + \alpha [R(t+1) + \text{discount} * \max Q(S_{t+1}, \underline{A}) - Q(S_t, A_t)]$

The difference lays in the lookup on Q' . For ***Q-Learning*** we have to find the maximum ***Q-value*** for the update equation by changing the \underline{A} , then we will have a new ***Q(S_t, A_t)***. This means we learn about the action-value function of the optimal policy, even when the behavior policy is not the optimal policy. For ***SARSA*** the action we learn is the actual action $\underline{A_{t+1}}$ that we follow.

In the end ***Q-Learning*** backs up the best Q-value from the state reached while ***SARSA*** waits until an action is taken and then backs up the Q-value from that action.