

Sarsa: $R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$ A_{t+1} selected
 target π must be b \downarrow action taken on-policy (e.g. ϵ -greedy)

Expected Sarsa target: $R_{t+1} + \gamma E_{\pi} [Q(S_{t+1}, A')]$
 $\sum_{a' \in A} \pi(a' | S_{t+1}) Q(S_{t+1}, a')$

π does not need to be b . \nwarrow does not include A_{t+1}

Q-learning: $R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a')$
 \uparrow does not include A_{t+1}