

SARSA vs Q learning

on vs off policy learning

on strategy: has two competing goals:

(1) exploration:

(2) exploitation

off-policy learning: use two strategies:

one for exploring and one for exploiting

We design a policy improvement algorithm in the statistical framework

Objective: learn q_{π} : $S \times A \rightarrow \mathbb{R}$

initialise q : $S \times A \rightarrow \mathbb{R}$ (0 or random)

REPEAT:

$S \leftarrow S_0$

while (trajectory not over):

choose action a using ϵ -greedy strategy
from q

sample $(s, a) : \mathcal{R}, s'$

CASE ①: SARSA (on-policy):

$(s, a, \mathcal{R}, s', a')$

evaluate using
 ϵ -greedy

$$q(s, a) = q(s, a) + \alpha [\mathcal{R} + \gamma q(s', a') - q(s, a)]$$

CASE ②: Q-learning (off-policy)

evaluating using greedy

$$q(s, a) = q(s, a) + \alpha [\mathcal{R} + \gamma \max_{a' \in A} q(s', a') - q(s, a)]$$

