

Q1. Yes. Proof:

$$\pi^* \stackrel{\text{by definition}}{=} \operatorname{argmax}_a Q_M^*(s, a) = \operatorname{argmax}_a Q_M^*(s, a) + \underbrace{f(s)}_{\text{const.}} = \operatorname{argmax}_a Q_{M'}^*(s, a)$$

$\pi^* = \operatorname{argmax}_a Q_{M'}^*(s, a)$ is defined as the ^{w.r.t. a} optimal policy

for $Q_M^*(s, a)$, so we proved that $\pi^*(s)$ is the optimal

Policy for $Q_{M'}^*(s, a)$ also!

Q3. 1. Each iteration we run π_t , the opt. policy for iter. t ,
until a new state-action is found. Let's assume by
contradiction that exists a pair (s', a') such that
 $R'_t(s', a') = 1$ we didn't discover at iteration t . In other words, the
optimal policy π_t stuck on $(s, a) : R'_t(s, a) = 0$, so $V_t^\pi = 0$.

Given a strongly connected MDP, we know (s', a') can be visited.
So a policy π' that does so will get a reward of 1 at (s', a') ,
thus π is not optimal, in contradiction!

Each iteration must end!

2. $T_{\text{iteration}} = O(|S|)$. For example, in the case that we should go over all the states that we have already visited in order to go from S_0 to (s', a')

3. every state s in the set S has a set of possible action $A(s)$, So in order to go over the entire

S , a space we need $\sum_{s \in S} |A(s)| = |S| \cdot |A|$ iterations

4. $T_{\text{exploitation}} = T_{\text{iteration}} \cdot \text{num-iterations} = O(|S|) \cdot |S| \cdot |A| =$
 $= O(|S|^2 \cdot |A|)$