

Q3. Suppose G_t is return at time t for some episode following first visit MC.

We essentially want $E[p_{t:T-1} G_t | S_t, A_t]$

$= q_\pi(S_t, A_t)$, where

$$p_{t:T-1} = \prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{b(A_k | S_k)}, \text{ where}$$

π is the target policy and b is the behaviour policy. Then,

$$q_\pi(s, a) = \frac{\sum_{t \in \mathcal{T}(s, a)} p_{t:T-1} G_t}{|\mathcal{T}(s, a)|} \quad \text{where}$$

$\mathcal{T}(s, a) =$ The set of all time steps where
 state s is visited and action a was
 taken (ie, $(S_t, A_t) = (s, a)$)

Q2: Here 'o' represents a state, '•' represents a state action pair and \square is the terminal state to estimate $q_{\pi}(s,a)$ we will have the following for an episode:

• \rightarrow • \rightarrow • \rightarrow ... \rightarrow \square

Q3: Suppose for a new episode our estimates of time to becomes very small. TD updates will give more weight to these new updates as compared to MC (that is still marking against the ~~un~~ updated q_t)

Q5. So it is clear that the first episode ends at the left terminal state (T). This is because following TD(0):

$$V(S_t) = V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

for $V(E)$ would have updated to

$$\begin{aligned} V(E) &= 0.5 + 0.1(1 + 0 - 0.5) \\ &= 0.55. \end{aligned}$$

had the episode ended on the right.

Instead we have

$$V(X) = 0.5 + 0.1(0 + 1 \times 0.5 - 0.5)$$

for $X \in \{B, C, D, E\}$

$$\begin{aligned} \text{and } V(A) &= 0.5 + 0.1(0 + 1 \times 0 - 0.5) \\ &= 0.45. \end{aligned}$$

Q6: Ex. 6.3: Yes. Suppose if our initial estimates were very off. We would have a significant RMSE. In such a case a larger α would quickly reduce RMSE before converging onto some value. At the same time, a smaller RMSE α would mean a longer time to converge, but the steady value (of RMSE) would also be smaller.

There is probably a value of α so that there is faster convergence and a lower RMSE value.

Ideally we would keep a small α and observe that RMSE becomes 0 after an infinite no. of episodes.

x.6.4. ~~Yes~~ No, this ~~would always occur~~ ^{is always due to initialization}. Consider for any state x :

$$V(x) = V(x) + \alpha (R_{t+1} + \gamma V(y) - V(x))$$

$$\text{Then, } V(x) = V(x)(1-\alpha) + \alpha (R_{t+1} + \gamma V(y))$$

$$\text{and in our situation: } V(x)(1-\alpha) + \alpha (R_{t+1} + V(y))$$

where y is the new state transition

Now suppose we had a reasonable estimate $V'(x) \approx V(x) + \delta$ and also some estimate $V'(y)$

Here for most cases ($E \rightarrow \text{Terminal}$) $R_{t+1} = 0$. Ignoring that,

$$V(x) = V'(x)(1-\alpha) + \alpha \cdot V'(y). \text{ For a high}$$

α , we would completely lose our estimate $V'(x)$ and instead obtain $V(x) \approx V'(y)$.

which would give a high RMSE. Now suppose our initialization was perfect.

Now suppose our initialization was perfect. i.e. $V(x) = V_n(x)$. Our updates would very quickly make values approach $V(E)$.

$$V(E) = V(E)_{init} + \alpha (1 - V(E)_{init}) \\ = V(E)(1 - \alpha), \text{ or}$$

$$V(D) = V(E)_{init} + \alpha (0 + \gamma(V(D)_{init}) - V(E)_{init}) \\ = V(E)_{init} (1 - \alpha) + \alpha (V(D)_{init})$$

for $\alpha \approx 1$ we give too much weight to our updates.

Supposing we had a bad initialization, (say, $V(x) = 0.9$ for $x \in \{A, B, C, D, E\}$)

RMSE after update for some $\alpha = 0.8$, would

$$\text{see } V(x)' = 0.9 \times 0.2 \text{ or } 0.9 \times 0.2 + 0.9 \times 0.8 \\ = 0.18 \text{ or } 0.9$$

$$\text{On average: } \frac{0.18 + 0.9}{2} = 0.54$$

compared for same $\alpha = 0.2$, the same update gives : 0.9×0.8 or 0.9

$$= \frac{0.72 + 0.9}{2}$$

on average $\rightarrow 0.61$

(When policy is arbitrary)

The decrease initially in RMSE is more but increases due to instability.

and policy

Ex. 6.5.

In either case, we must know the ~~policy~~ ^{model} π .

By the copies of this it is likely the case that Policy ^{iteration} ~~evaluation~~ was done and $V_{\pi}(\cdot)$ values were obtained. The model can be inferred through and we can calculate

$$V_{\pi}(s) = E_{\pi}(R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s)$$

The other way to do this would be via Monte Carlo methods for an infinite no. of episodes.

Ex. 6.12: No they are still not the same. Even with greedy action selection, SARSA remains on-policy while Q-learning is off-policy.

Also, as per the pseudo-code SARSA updates Q after choosing S', A' , while Q-learning is first updating Q , then choosing the next action.