

Q.1 a) in online Q learning, all $Q(s,a)$ pairs are parameterized by same weight parameters θ .
So upon updating θ , will change $Q(s,a)$ of all (s,a) pairs and won't retain the old values.

In Watkins Q learning, all (s,a) pairs have their separate $Q(s,a)$ values. So updating value of one (s,a) pair doesn't affect the Q value of another pair (s',a') .

b) ~~Again, the update will suffer~~
Since the replay buffer here is static. The stochastic batch update will make the approximator Q_ϕ to only learn Q values corresponding to the replay buffer we have.
As Q_ϕ keeps updating, we need to collect more replay samples using current Q function to help the agent learn better and converge to the optimal state-value function.

Since after each update, policy changes, one need to collect samples from updated policy not old one.

Also, if target networks are not used, the networks will further face convergence issues.

$$Q.3(d) \quad \nabla_{\theta} J(\theta) = E_{\pi_{\theta}} \left[\sum_{t=0}^{\infty} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \psi_t \right]$$

$$\psi_t = G_{t:\infty} - b(s_t)$$

$$\nabla_{\theta} J(\theta) = E_{\pi_{\theta}} \left[\sum_{t=0}^{\infty} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) (G_{t:\infty} - b(s_t)) \right]$$

writing in terms of trajectories τ

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[\nabla_{\theta} \log \pi_{\theta}(\tau) (G_{\tau:\infty} - b) \right]$$

$$\text{Variance (Var)} = E_{\tau \sim \pi_{\theta}(\tau)} \left[\left(\nabla_{\theta} \log \pi_{\theta}(\tau) (G_{\tau:\infty} - b) \right)^2 \right] - \underbrace{E_{\tau \sim \pi_{\theta}(\tau)} \left[\nabla_{\theta} \log \pi_{\theta}(\tau) (G_{\tau:\infty} - b) \right]}_{\text{baselines are unbiased}}$$

so it is just $E_{\tau \sim \pi_{\theta}(\tau)} \left[\nabla_{\theta} \log \pi_{\theta}(\tau) G_{\tau:\infty} \right]$

$$\text{Var} = E_{\tau \sim \pi_{\theta}(\tau)} \left[\left(\nabla_{\theta} \log \pi_{\theta}(\tau) (G_{\tau:\infty} - b) \right)^2 \right] - E_{\tau \sim \pi_{\theta}(\tau)} \left[\nabla_{\theta} \log \pi_{\theta}(\tau) G_{\tau:\infty} \right]$$

$$\frac{d(\text{Var})}{db} = -2 E_{\tau \sim \pi_{\theta}(\tau)} \left[\nabla_{\theta} \log \pi_{\theta}(\tau) G_{\tau:\infty} \right] + 2b E_{\tau \sim \pi_{\theta}(\tau)} \left[\nabla_{\theta} \log \pi_{\theta}(\tau) \right]$$

$$\frac{d(\text{Var})}{db} = 0 \Rightarrow b = \frac{E_{\tau \sim \pi_{\theta}(\tau)} \left[\nabla_{\theta} \log \pi_{\theta}(\tau) G_{\tau:\infty} \right]}{E_{\tau \sim \pi_{\theta}(\tau)} \left[\nabla_{\theta} \log \pi_{\theta}(\tau) \right]}$$