« Alireza Gargoori Motlegh _ 98102176 »

« Homework 3 _ Reinforcement Learning »

**1**

**(a)** One of the main requirements of $LQR$ to converge is that environment must be fully-observable and we have $x_t$ in each time step. Also, the quadratic approximation for the cost function and linear approximation for dynamics must be reasonable. Another thing is that in short time horizons $LQR$ may lead to suboptimal solutions. Also, the system needs to be controllable in full coordinates.

**(b)** When the environment is partially observable, we have POMDP and therefore we can estimate states through the observations using HMMs : $P(s_t | O_{1:t}) = P(O_t | s_t) \sum_{s_{t-1}} P(s_t | s_{t-1}) P(s_{t-1} | O_{1:t-1})$

or Variational inference : treating hidden states as latent variables.

After estimating the states, we can replace them and run LQR; it can be shown that optimal policy is : $u_t = K_t \mathbb{E}[x_t | O_{1:t}] + k_t$

(c) Model-free methods suffer from sample-efficiency while model-based methods suffer from significant bias, since complex unknown dynamics cannot always be modeled accurately enough to produce effective; therefore combining these methods could result in a better unbiased sample-efficient model.
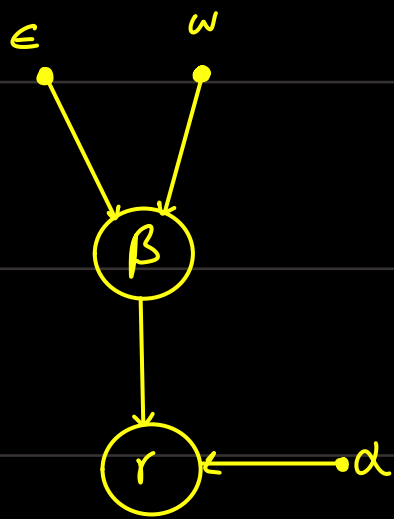
To combine LQR with a deep model-free model, we can ① use the model-free method (such as TPO, TRPO, SAC, DQN, etc.) to learn an initial control policy through interacting with environment ② ; then this learned policy could be used to collect trajectories for fitting the system's dynamics. ③ LQR assumes a linear dynamic for the system; so using the collected data, we fit a linear model to state-action to model the dynamics. ④ After designing a reasonable cost function and a quadratic approximation for that, ⑤ we can now run LQR for trajectory optimization and finding the optimal control policy. ⑥ Then we use the LQR solution by using it as a starting policy to fine-tune model-free policy and iterate the above steps (2-6) to improve the policy.

(d) In the stochastic system, we have $x_{t+1} \sim \mathbb{P}(x_{t+1} \mid x_t, u_t)$; e.g. a normal distribution; therfore we can model the uncertainty by using stochastic dynamics. However, in iLQR we can have a nonlinear estimation for the mean & std of this distribution and then approximate it as a Linear-quadratic for dynamics and Cost using taylor expansion around the current state-action pair. Other nonlinear parameterization of other distributions are possible, yet very hard to solve.

To ensure exploration, we can add a regularization term to the Cost function which accounts for exploration, such as entropy or information gain. These terms could be expanded through quadratic approximation as well and the stochastic iLQR would both capture the optimal control policy and enough exploration.

(a) □ $r \sim Gamma(\alpha, \beta) \longrightarrow \mathbb{P}(r|\alpha, \beta) = \frac{\beta^{\alpha}}{\Gamma(\alpha)} r^{\alpha-1} e^{-\beta r}$

$\beta \sim Gamma(\epsilon, \omega) \longrightarrow \mathbb{P}(\beta|\epsilon, \omega) = \frac{\omega^{\epsilon}}{\Gamma(\epsilon)} \beta^{\epsilon-1} e^{-\omega\beta}$

$$\mathbb{P}(\beta|r_1, \alpha, \epsilon, \omega) \propto \mathbb{P}(r_1|\alpha, \beta) \, \mathbb{P}(\beta|\epsilon, \omega)$$

$$= \frac{\beta^{\alpha}}{\Gamma(\alpha)} r^{\alpha-1} e^{-\beta r} \frac{\omega^{\epsilon}}{\Gamma(\epsilon)} \beta^{\epsilon-1} e^{-\omega\beta} = \frac{f(r_1, \omega, \alpha, \epsilon)}{\Gamma(\alpha+\epsilon)} \beta^{\alpha+\epsilon-1} e^{-(r_1+\omega)\beta}$$
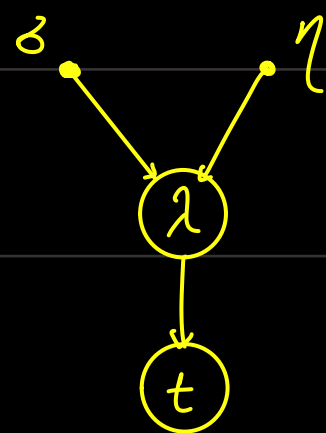
$\longrightarrow$ We can confirm that the $\mathbb{P}(\beta|\epsilon, \omega)$ is a conjugate prior

for the likelihood and therefore:

$$\mathbb{P}(\beta|r_1, \alpha, \epsilon, \omega) = \frac{(\omega')^{\epsilon'}}{\Gamma(\alpha+\epsilon)} \beta^{\overbrace{\alpha+\epsilon}^{\epsilon'}-1} e^{-\overbrace{(r_1+\omega)}^{\omega'}\beta} = \frac{(r_1+\omega)^{\alpha+\epsilon}}{\Gamma(\alpha+\epsilon)} \beta^{\alpha+\epsilon-1} e^{-(r_1+\omega)\beta}$$

$\boxed{\longrightarrow \quad \beta|r_1, \alpha, \epsilon, \omega \sim Gamma\left(\epsilon' = \alpha+\epsilon, \; \omega' = r_1+\omega\right)}$

□ $t \sim \text{Exp}(\lambda) \implies p(t|\lambda) = \lambda e^{-\lambda t}$

$\lambda \sim \text{Gamma}(\delta, \eta) \implies p(\lambda|\delta, \eta) = \dfrac{\eta^{\delta}}{\Gamma(\delta)} \lambda^{\delta-1} e^{-\eta \lambda}$



$$p(\lambda|t_i, \delta, \eta) \propto p(t_i|\lambda) \, p(\lambda|\delta, \eta)$$

$$= \lambda e^{-\lambda t_i} \frac{\eta^{\delta}}{\Gamma(\delta)} \lambda^{\delta-1} e^{-\eta \lambda} = f(\eta, \delta, t_i) \, \lambda^{\delta} \, e^{-(t_i+\eta)\lambda}$$

$\longrightarrow$ We can confirm that the $p(\lambda|t, \delta, \eta)$ is a conjugate prior

for the likelihood and therefore:

$$p(\lambda|t_i, \delta, \eta) = \frac{(\eta')^{\delta'}}{\Gamma(\delta')} \lambda^{\delta'-1} e^{-\eta' \lambda} \begin{cases} \delta' = \delta + 1 \\ \eta' = t_i + \eta \end{cases}$$

$\longrightarrow$ $\boxed{\begin{array}{l} \lambda|t_i, \delta, \eta \sim \text{Gamma}(\delta' = \delta+1, \, \eta' = t_i + \eta) \\[2em] p(\lambda|t_i, \delta, \eta) = \dfrac{(t_i+\eta)^{\delta+1}}{\Gamma(\delta+1)} \lambda^{\delta} e^{-(t_i+\eta)\lambda} \end{array}}$

(b) $\quad P(t_2|t_1) = \int P(t_2, \lambda | t_1)\, d\lambda = \int P(t_2|\lambda, t_1)\, P(\lambda|t_1)\, d\lambda =$

$$\int P(t_2|\lambda)\, P(\lambda|t_1)\, d\lambda = \int \lambda e^{-\lambda t_2}\, \frac{\eta'^{\delta'}}{\Gamma(\delta')}\, \lambda^{\delta'-1}\, e^{-\eta'\lambda}\, d\lambda =$$

$$\frac{\eta'^{\delta'}}{\Gamma(\delta')} \int_0^\infty \lambda^{\delta'}\, e^{-(\eta'+t_2)\lambda}\, d\lambda = \qquad \left(\Gamma(\delta) = (\delta-1)!\right)$$
$$\delta' \in \mathbb{N}$$

$$\frac{\eta'^{\delta'}}{(\delta-1)!}\, \frac{1}{(\eta'+t_2)} \underbrace{\int_0^\infty \lambda^{\delta'}\, (\eta'+t_2)\, e^{-(\eta'+t_2)\lambda}\, d\lambda}_{I} =$$

$I$ : Integration by parts $\begin{cases} u(\lambda) := \lambda^{\delta'} \\ v'(\lambda) = e^{-(\eta'+t_2)\lambda} \end{cases}$

$$I = \int_0^\infty u(\lambda)\, v'(\lambda)\, d\lambda = u(\lambda) v(\lambda)\Big|_0^\infty - \int_0^\infty u'(\lambda)\, v(\lambda)\, d\lambda$$

$$= \underbrace{\lambda^{\delta'}\, e^{-(\eta'+t_2)\lambda}\Big|_0^\infty}_{= 0 \ \left(\begin{array}{l}\text{since } \delta' \in \mathbb{N}, \\ \eta', t_2 > 0\end{array}\right)} - \frac{\delta'}{\eta'+t_2} \underbrace{\int_0^\infty \lambda^{\delta'-1}\, (\eta'+t_2)\, e^{-(\eta'+t_2)\lambda}\, d\lambda}_{II}$$

$$\text{II} = \underbrace{\lambda^{\delta'-1} e^{-(\eta'+t_2)\lambda} \Big|_0^\infty}_{= 0} - \frac{\delta'-1}{\eta'+t_2} \underbrace{\int_0^\infty \lambda^{\delta'-2} (\eta'+t_2) e^{-(\eta+t_2)\lambda} \, d\lambda}_{\text{III}}$$

$$\text{III} = \underbrace{\lambda^{\delta'-2} e^{-(\eta'+t_2)\lambda} \Big|_0^\infty}_{= 0} - \frac{\delta'-2}{\eta'+t_2} \underbrace{\int_0^\infty \lambda^{\delta'-2} (\eta'+t_2) e^{-(\eta+t_2)\lambda} \, d\lambda}_{\text{IV}}$$

$$\vdots$$

$$\delta' \text{ times} = \underbrace{\lambda^0 e^{-(\eta'+t_2)} \Big|_0^\infty}_{= 0} - \frac{1}{\eta'+t_2} \int_0^\infty (\eta'+t_2) e^{-(\eta'+t_2)\lambda} \, d\lambda$$

$$= \frac{-1}{\eta'+t_2} e^{-(\eta'+t_2)\lambda} \Big|_0^\infty = \frac{1}{\eta'+t_2}$$

$$\Rightarrow \boxed{\; \text{I} = \frac{\delta'}{\eta'+t_2} \left( \frac{\delta'-1}{\eta'+t_2} \left( \frac{\delta'-2}{\eta'+t_2} \left( \cdots \frac{1}{\eta'+t_2} \right) \right) \right) = \frac{\delta'!}{(\eta'+t_2)^{\delta'}} \;}$$

$$\Rightarrow P(t_2 \mid t_1) = \frac{\eta'^{\delta'}}{(\delta-1)!} \frac{1}{(\eta'+t_2)} \frac{\delta!}{(\eta'+t_2)^{\delta'}} = \frac{\eta'^{\delta'} \delta'}{(\eta'+t_2)^{\delta'+1}}$$

$$\longrightarrow P(t_2 \mid t_1) = \frac{\eta'^{\delta'} \delta' \times \eta'^{-(\delta'+1)}}{(\eta' + t_2)^{\delta'+1} \times \eta'^{-(\delta'+1)}} = \frac{\eta'^{-1} \delta'}{\left(1 + \frac{t_2}{\eta'}\right)^{\delta'+1}}$$

$$= \frac{\delta'}{\eta'} \left(1 + \frac{t_2}{\eta'}\right)^{-(\delta'+1)}$$

$$\Longrightarrow \boxed{P(t_2 \mid \delta', \eta') = \frac{\delta'}{\eta'} \left(1 + \frac{t_2}{\eta'}\right)^{-(\delta'+1)}}$$

$$\boxed{t_2 \sim \text{Lomax}(\delta', \eta')}$$

(C) <u>Risky</u>: Since the rewards could get to infinity, in this scenario the agent would always choose to play. <u>Risk-free</u>: In contrast to the previous, since the time it takes could be anything from zero to infinity, the agent would always choose to not play.

<u>Risk-neutral</u> : we need to compare $\mathbb{E}\left[\frac{r}{t}\right]$ with $K$ ;

$$\mathbb{E}\left[r/t\right] = \mathbb{E}\left[r\right]\mathbb{E}\left[\frac{1}{t}\right]$$

$$\mathbb{E}\left[r\right] = \int r \int P(r|\alpha,\beta) P(\beta|\epsilon,\omega) \, d\beta \, dr = R$$

$$\mathbb{E}\left[\frac{1}{t}\right] = \int \frac{1}{t} \int P(t|\lambda) P(\lambda|\sigma,\eta) \, d\lambda \, dt = T$$

if $\quad \frac{R}{T} \geqslant K \implies$ would play

else $\qquad \implies$ wouldn't play .

(a)

$$\log P(O_{1:T}) \geqslant \sum_t \mathbb{E}_{(s_t, a_t) \sim q} \left[ r(s_t, a_t) + H(q(a_t | s_t)) \right]$$

$t = T:$

$$\mathbb{E}_{(s_T, a_T) \sim q} \left[ r(s_T, a_T) - \log q(a_T | s_T) \right] =$$

$$\mathbb{E}_{(s_T, a_T) \sim q} \left[ \log(\exp(r(s_T, a_T)) - \log q(a_T | s_T) \right] =$$

$$\mathbb{E}_{s_T \sim q(s_T)} \left[ \mathbb{E}_{a_T \sim q(a_T | s_T)} \left[ \log(\exp(r(s_T, a_T)) - \log q(a_T | s_T) \right] \right.$$

$$= \mathbb{E}_{s_T \sim q(s_T)} \left[ -D_{KL}\left( q(a_T | s_T) \,\|\, \frac{1}{Z} \exp(r(s_T, a_T)) \right) \right]$$

We wanted to maximize lower bound; since

$D_{KL} \geqslant 0$, we need to have:

$$q(a_T | s_T) = \arg\max_q -D_{KL}\left( q(a_T | s_T) \,\|\, \frac{1}{Z} \exp(r(s_T, a_T)) \right)$$

$$= \arg\min_q D_{KL}\left( q(a_T | s_T) \,\|\, \frac{1}{Z} \exp(r(s_T, a_T)) \right) \Rightarrow$$

$$q(a_t | s_t) = \frac{1}{Z} \exp\left(r(s_T, a_T)\right)$$

$$Z = \int \exp\left(r(s_T, a_T)\right) da_T$$

$$\Rightarrow q(a_t | s_t) = \frac{\exp\left(r(s_T, a_T)\right)}{\int \exp(r(s_T, a_T)) da_T}$$

for $t = T$ we have: $Q(s_T, a_T) = r(s_T, a_T)$ and

therefore:

$$V(s_T) = \log \int \exp\left(Q(s_T, a_T)\right) da_T = \log \int \exp(r(s_T, a_T)) da_T$$

$$\boxed{q(a_T | s_T) = \frac{\exp\left(Q(s_T, a_T)\right)}{\exp\left(V(s_T)\right)} = \exp\left(Q(s_T, a_T) - V(s_T)\right)}$$

for $t < T$:

$$q(a_t | s_t) = \arg\max_q \mathbb{E}_{s_t \sim q(s_t)}\left[ \underbrace{\mathbb{E}_{a_t \sim q(a_t | s_t)}\left[ r(s_t, a_t) + \mathbb{E}_{P(s_{t+1} | a_t, s_t)}\left[ V(s_{t+1}) \right]\right]}_{Q(s_t, a_t)} - \log q(a_T | s_T)\right]$$

$$= \ldots \text{ (next page)}$$

$$\longrightarrow \quad q(a_t|s_t) = \underset{q}{\arg\max} \; \underset{s_t \sim q(s_t)}{\mathbb{E}}\left[ \underset{a_t \sim q(a_t|s_t)}{\mathbb{E}}\left[ Q(s_t, a_t) - \log q(a_t|s_t) \right]\right]$$

$$= \underset{q}{\arg\min} \; \underset{s_t \sim q(s_t)}{\mathbb{E}}\left[ D_{KL}\left( q(a_t|s_t) \;||\; \frac{1}{Z} \exp\left( Q(s_t, a_t) \right) \right) \right]$$

$$D_{KL} \geqslant 0 \quad \text{and} \quad D_{KL}(p||q) = 0 \quad \text{when } p=q \implies$$

$$q(a_t|s_t) = \frac{1}{Z} \exp\left( Q(s_t, a_t) \right) \; , \quad Z = \int \exp\left( Q(s_t, a_t) \right) da_t$$

$$\implies \quad q(a_t|s_t) = \frac{\exp\left( Q(s_t, a_t) \right)}{\int \exp\left( Q(s_t, a_t) \right) da_t}$$

$$V(s_t) = \log \int \exp\left( Q(s_t, a_t) \right) da_t \implies$$

$$q(a_t|s_t) = \frac{\exp\left( Q(s_t, a_t) \right)}{\exp\left( V(s_t) \right)} = \exp\left( Q(s_t, a_t) - V(s_t) \right)$$

$$\boxed{\implies \quad q(a_t|s_t) = \exp\left( Q(s_t, a_t) - V(s_t) \right) = \exp\left( A(s_t, a_t) \right)}$$

(b)

$$\pi_\theta(s_t, a_t) = \pi_\theta(a_t \mid s_t)\, \pi_\theta(s_t)$$

$$J(\theta) = \sum_t \mathbb{E}_{(s_t, a_t) \sim \pi_\theta(s_t, a_t)} \left[ r(s_t, a_t) + H(\pi_\theta(a_t \mid s_t)) \right]$$

$$\Rightarrow \nabla_\theta J(\theta) = \nabla_\theta \sum_t \mathbb{E}_{(s_t, a_t) \sim \pi_\theta(s_t, a_t)} \left[ r(s_t, a_t) + H(\pi_\theta(a_t \mid s_t)) \right]$$

$$= \sum_t \nabla_\theta \mathbb{E}_{(s_t, a_t) \sim \pi_\theta(s_t, a_t)} \left[ r(s_t, a_t) - \log(\pi_\theta(a_t \mid s_t)) \right]$$

$$\Rightarrow \nabla_\theta J(\theta) = \nabla_\theta \sum_\tau \pi_\theta(\tau) \sum_t \left( r(s_t, a_t) - \log \pi_\theta(a_t \mid s_t) \right)$$

$$= \underbrace{\nabla_\theta \sum_\tau \pi_\theta(\tau) \sum_t r(s_t, a_t)}_{\text{I}} - \underbrace{\nabla_\theta \sum_\tau \pi_\theta(\tau) \sum_t \log \pi_\theta(a_t \mid s_t)}_{\text{II}}$$

I:
$$\nabla_\theta \sum_\tau \pi_\theta(\tau)\, r(\tau) = \sum_\tau \nabla_\theta \pi_\theta(\tau)\, r(\tau) = \sum_\tau \pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau)\, r(\tau)$$

$$= \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \nabla_\theta \log \pi_\theta(\tau) \underbrace{\sum_t r(s_t, a_t)}_{*} \right]$$

$*$
$$\pi_\theta(\tau) = p(s_1) \prod_t p(s_{t+1} \mid s_t, a_t)\, \pi_\theta(a_t \mid s_t) \Rightarrow \nabla_\theta \log \pi_\theta(\tau) = \sum_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

$$\Rightarrow \boxed{ \text{I} = \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \sum_t \nabla_\theta \log \pi_\theta(a_t \mid s_t) \sum_t r(s_t, a_t) \right] }$$

$\mathrm{II}:$

$$\nabla_\theta \sum_\tau \pi_\theta(\tau) \sum_t \log \pi_\theta(a_t|s_t) =$$

$$\sum_\tau \nabla_\theta \pi_\theta(\tau) \sum_t \log \pi_\theta(a_t|s_t) + \sum_\tau \pi_\theta(\tau) \underbrace{\sum_t \nabla_\theta \log \pi_\theta(a_t|s_t)}_{= \nabla_\theta \log \pi_\theta(\tau)} =$$

$$\sum_\tau \nabla_\theta \pi_\theta(\tau) \sum_t \log \pi_\theta(a_t|s_t) + \underbrace{\sum_\tau \pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau)}_{\nabla_\theta \pi_\theta(\tau)} =$$

$$\sum_\tau \nabla_\theta \pi_\theta(\tau) \left( \sum_t \log \pi_\theta(a_t|s_t) + 1 \right) =$$

$$\sum_\tau \pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau) \left( \sum_t \log \pi_\theta(a_t|s_t) + 1 \right) =$$

$$\mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \nabla_\theta \log \pi_\theta(\tau) \left( \sum_t \log \pi_\theta(a_t|s_t) + 1 \right) \right]$$

$\ast$

$\Longrightarrow$

$$\boxed{\mathrm{II} = \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \sum_t \nabla_\theta \log \pi_\theta(a_t|s_t) \left( \sum_t \log \pi_\theta(a_t|s_t) + 1 \right) \right]}$$

$$\nabla_\theta J(\theta) = \mathrm{I} - \mathrm{II} \quad \Longrightarrow$$

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \sum_t \nabla_\theta \log \pi_\theta(a_t | s_t) \sum_t r(s_t, a_t) \right] -$$

$$\mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \sum_t \nabla_\theta \log \pi_\theta(a_t | s_t) \sum_t \log \pi_\theta(a_t | s_t) + 1 \right] =$$

$$\mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(a_t | s_t) \left( \sum_{t=1}^{T} \left[ r(s_t, a_t) - \log \pi_\theta(a_t | s_t) \right] - 1 \right) \right]$$

Due to causality, that is the policy at $\underline{t}$ cannot affect the previous rewards and policies, we can rewrite the equation as:

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(a_t | s_t) \left( \sum_{t'=t}^{T} \left[ r(s_{t'}, a_{t'}) - \log \pi_\theta(a_{t'} | s_{t'}) \right] - 1 \right) \right]$$

$$= \mathbb{E}_{(s_t, a_t) \sim \pi_\theta(s_t, a_t)} \left[ \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(a_t | s_t) \left( \sum_{t'=t}^{T} \left[ r(s_{t'}, a_{t'}) - \log \pi_\theta(a_{t'} | s_{t'}) \right] - 1 \right) \right]$$

$\implies$ We can sample $N$ trajectories and approximate $\nabla_\theta J(\theta)$ as:

$$\nabla_\theta J(\theta) \simeq \frac{1}{N} \sum_i \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(a_{it} | s_{it}) \left( \sum_{t'=t}^{T} \left[ r(s_{i,t'}, a_{i,t'}) - \log \pi_\theta(a_{i,t'} | s_{i,t'}) \right] - 1 \right)$$

**(c)**

$$\pi_\theta(a_t | s_t) = \exp\left( Q_{(s_t, a_t)} - V_{(s_t)} \right) \quad *$$

$$\nabla_\theta \mathcal{J}(\theta) \simeq \frac{1}{N} \sum_i \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \left( \sum_{t'=t}^T \left[ r(s_{t'}, a_{t'}) - \log \pi_\theta(a_t | s_t) \right] - 1 \right) \quad **$$

$$\nabla_\theta J(\theta) \simeq \frac{1}{N} \sum_i \sum_{t=1}^T \nabla_\theta \left( Q_{(s_t, a_t)} - V_{(s_t)} \right) \left( \sum_{t'=t}^T \left[ r(s_{t'}, a_{t'}) - Q_{(s_{t'}, a_{t'})} + V_{(s_{t'})} \right] - 1 \right)$$

$$= \frac{1}{N} \sum_i \sum_{t=1}^T \left( \nabla_\theta Q_{(s_t, a_t)} - \nabla_\theta V_{(s_t)} \right) \left( \sum_{t'=t}^T \left[ r(s_{t'}, a_{t'}) - Q_{(s_{t'}, a_{t'})} + V_{(s_{t'})} \right] - 1 \right)$$

$$\boxed{ \nabla_\theta J(\theta) \simeq \frac{1}{N} \sum_i \sum_{t=1}^T \left( \nabla_\theta Q_{(s_t, a_t)} - \nabla_\theta V_{(s_t)} \right) \left( \underbrace{\sum_{t'=t}^T \left[ r(s_{t'}, a_{t'}) - Q_{(s_{t'}, a_{t'})} + V_{(s_{t'})} \right] - 1 }_{ -\mathbb{E}[V_{(s_{t+1})}]} \right) }$$

**(d)**

$*$ MaxEnt RL grad. :

$$\boxed{ \begin{aligned} \nabla_\theta J(\theta) &= \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \sum_{t=1}^T \left( \nabla_\theta Q_{(s_t, a_t)} - \nabla_\theta V_{(s_t)} \right) \hat{A}_{(s_t, a_t)} \right] \\ &\simeq \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \left( \nabla_\theta Q_{(s_t^{(i)}, a_t^{(i)})} - \nabla_\theta V_{(s_t^{(i)})} \right) \hat{A}_{(s_t^{(i)}, a_t^{(i)})} \end{aligned} }$$

$*$ soft Q-learning grad. :

$$\boxed{ \begin{aligned} \nabla_\theta J(\theta) &= \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \sum_{t=1}^T \nabla_\theta Q_{(s_t, a_t)} \hat{A}_{(s_t, a_t)} \right] \\ &\simeq \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta Q_{(s_t^{(i)}, a_t^{(i)})} \hat{A}_{(s_t^{(i)}, a_t^{(i)})} \end{aligned} }$$

There is an extra term, $-\nabla_\theta V(s_t)$ in MaxEnt RL gradient which resemble the insufficiency of policy gradient to resolve the addition or substraction of an action-independent constant.

This term could be eliminated for a particular choice of $\hat{A}(s_t, a_t)$, which in the paper "equivalence between policy gradients and soft Q-learning" is explained in detail.

In general, Q-learning methods are more sample-efficient but their Q-function estimate may be inaccurate; but in the sense of entropy-regularized Q-learning has a better estimate of this function.