

Naser Kazemi 99102059 HW3

(Q1)

A) There are two possible sources for the randomness in the cumulative reward:

- 1) The policy randomness. If the policy itself is stochastic, then we have to compute the expected value.
- 2) The environment dynamics. The transition between states after taking an action can be probabilistic.

B) One way could be reparametrizing the actions. i.e., let a be an action sampled from the policy $\pi_\theta(a|s)$. We can reparametrize a as $a = f_\theta(s, \epsilon)$ where $\epsilon \sim p(\epsilon)$. Let's assume p is a gaussian distribution, with mean $\mu_\theta(s)$ and variance $\sigma_\theta^2(s)$, then

$$a = f_\theta(s, \epsilon) = \mu_\theta(s) + \sigma_\theta(s)\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$

Now we can see the dependency of rewards and θ , and we have

$$J_\theta = E \left[\sum_t \gamma^t r(s_t, f_\theta(s_t, \epsilon)) \right]$$

which is the expectation of a function of θ , and hence we can apply gradient ascent directly on the objective function.

(Q2)

A) I consider the following assumptions in my proof:

- 1) Trajectories are independent
- 2) Rewards and state-actions, i.e. $r(s_t, a_t)$ are consistent
- 2) The policy $\pi_\theta(a|s)$ is differentiable.

We have

$$\begin{aligned}
 D_\theta J(\theta) &= D_\theta E_{\tau \sim P(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \\
 &= D_\theta \int P_\theta(\tau) R(\tau) d\tau & R(\tau) = \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \\
 &= \int D_\theta P_\theta(\tau) R(\tau) d\tau & \tau = (s_0, a_0, s_1, a_1, \dots) \\
 &= \int \frac{D_\theta P_\theta(\tau)}{P_\theta(\tau)} P_\theta(\tau) R(\tau) d\tau \\
 &= E_{\tau \sim P_\theta(\tau)} [D_\theta \log P_\theta(\tau) R(\tau)] \\
 &= E_{\tau \sim P_\theta(\tau)} [\sum_{t=0}^{\infty} D_\theta \log \pi_\theta(s_t, a_t)]
 \end{aligned}$$

$$P_\theta(\tau) = \prod_{t=0}^{\infty} \pi_\theta(a_t | s_t) P(s_{t+1} | a_t, s_t) \Rightarrow D_\theta \log P_\theta(\tau) = \sum_{t=0}^{\infty} D_\theta \log \pi_\theta(a_t | s_t)$$

Now on the other hand:

$$\begin{aligned}
 &E_{\tau \sim P_\theta(\tau)} \left[\frac{1}{N} \sum_{i=1}^N \left(\sum_{t=0}^T D_\theta \log \pi_\theta(a_{i,t} | s_{i,t}) \right) \left(\sum_{t=0}^T \gamma^t r(s_{i,t}, a_{i,t}) \right) \right] \\
 &= \frac{1}{N} \sum_{i=1}^N E_{\tau \sim P_\theta(\tau)} \left[\left(\sum_{t=0}^T D_\theta \log \pi_\theta(a_{i,t} | s_{i,t}) \right) \left(\sum_{t=0}^T \gamma^t r(s_{i,t}, a_{i,t}) \right) \right]
 \end{aligned}$$

$$= \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} J(\theta) \\ = \nabla_{\theta} J(\theta)$$

In this proof, we also assume that the maximum length of each trajectory is T .

B) First let's prove a lemma:

EGLP Lemma: Suppose P_{θ} is a parametrized distribution over a random variable x . Then

$$E_{x \sim P_{\theta}} [\nabla_{\theta} \log P_{\theta}(x)] = 0$$

Proof: We know that

$$\int_x P_{\theta}(x) dx = 1$$

and

$$0 = \nabla_{\theta} 1 = \nabla_{\theta} \int_x P_{\theta}(x) dx \\ = \int_x \nabla_{\theta} P_{\theta}(x) dx \\ = \int_x \frac{\nabla P_{\theta}(x)}{P_{\theta}(x)} P_{\theta}(x) dx \\ = E_{x \sim P_{\theta}} [\nabla \log P_{\theta}(x)]$$

Now let's get back to the main problem.

$$E_T \left[\left(\sum_{t=0}^{\infty} D_\theta \log \pi_\theta(a_t | s_t) \right) \left(\sum_{T=0}^{\infty} \gamma^t r(s_t, a_t) \right) \right] \\ = \sum_{t=0}^{\infty} E_T \left[D_\theta \log \pi_\theta(a_t | s_t) \left(\sum_{t'=0}^{\infty} \gamma^{t'} r(s_{t'}, a_{t'}) \right) \right]$$

We have:

$$E_T \left[D_\theta \log \pi_\theta(a_t | s_t) \left(\sum_{t'=0}^{\infty} \gamma^{t'} r(s_{t'}, a_{t'}) \right) \right] \\ = E_T \left[D_\theta \log \pi_\theta(a_t | s_t) \sum_{t'=0}^{t-1} \gamma^{t'} r(s_{t'}, a_{t'}) \right] \\ + E_T \left[D_\theta \log \pi_\theta(a_t | s_t) \sum_{t'=t}^{\infty} \gamma^{t'} r(s_{t'}, a_{t'}) \right]$$

The random variables $\sum_{t'=0}^{t-1} \gamma^{t'} r(s_{t'}, a_{t'})$ and $D_\theta \log \pi_\theta(a_t | s_t)$ are independent because of the Markov property, therefore we can write

$$E_T \left[D_\theta \log \pi_\theta(a_t | s_t) \sum_{t'=0}^{t-1} \gamma^{t'} r(s_{t'}, a_{t'}) \right] \\ = E_{s_0, a_0, \dots, s_{t-1}, a_{t-1}} \left[\sum_{t'=0}^{t-1} \gamma^{t'} r(s_{t'}, a_{t'}) \right] E_{s_t, a_t, \dots} [D_\theta \log \pi_\theta(a_t | s_t)] \\ = 0$$

since by EGLP, $E_{s_t, a_t, \dots} [D_\theta \log \pi_\theta(a_t | s_t)] = 0$.

Thus,

$$E_T \left[D_\theta \log \pi_\theta(a_t | s_t) \left(\sum_{t'=t}^{\infty} \gamma^{t'} r(s_{t'}, a_{t'}) \right) \right] \\ = \gamma^t E_T \left[D_\theta \log \pi_\theta(a_t | s_t) \sum_{t'=t}^{\infty} \gamma^{t'-t} r(s_{t'}, a_{t'}) \right]$$

and

$$\begin{aligned}
& \mathbb{E}_T \left[\left(\sum_{t=0}^{\infty} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left(\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right) \right] \\
&= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_T \left[\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t'=t}^{\infty} \gamma^{t'-t} r(s_{t'}, a_{t'}) \right] \\
&= \sum_{t=0}^{\infty} \gamma^t \int \frac{\nabla_{\theta} \pi_{\theta}(a_t | s_t)}{\pi_{\theta}(a_t | s_t)} \prod_{t'=0}^{\infty} \pi_{\theta}(a_{t'} | s_{t'}) p(s_{t+1} | s_t, a_t) \\
&\quad \left(\sum_{t'=t}^{\infty} \gamma^{t'-t} r(s_{t'}, a_{t'}) \right) dt \\
&= \sum_{t=0}^{\infty} \gamma^t \int \frac{\nabla_{\theta} \pi_{\theta}(a_t | s_t)}{\pi_{\theta}(a_t | s_t)} \prod_{t'=0}^{\infty} \pi_{\theta}(a_{t'} | s_{t'}) p(s_{t+1} | s_t, a_t) Q(s_t, a_t) dt \\
&= \sum_{t=0}^{\infty} \gamma^t \int_S \nabla_{\theta} \log \pi_{\theta}(a_t = a | s_t = s) P(s_t = s) \pi_{\theta}(a_t = a | s_t = s) Q(s_t = s, a_t = a) ds da \\
&= \sum_{t=0}^{\infty} \gamma^t \int_S P(s_t = s) ds \int_a \nabla_{\theta} \log \pi_{\theta}(a_t = a | s_t = s) \pi_{\theta}(a_t = a | s_t = s) Q(s_t = s, a_t = a) da \\
&= \int_S \sum_{t=0}^{\infty} \gamma^t P(s_t = s) ds \int_a \nabla_{\theta} \log \pi_{\theta}(a_t = a | s_t = s) \pi_{\theta}(a_t = a | s_t = s) Q(s_t = s, a_t = a) da \\
&= \int_S P_{\pi}(s) \mathbb{E}_{\pi} [\nabla_{\theta} \log \pi_{\theta}(a | s) Q(s, a)] \\
&= \mathbb{E}_{P_{\pi} \otimes \pi} [\nabla_{\theta} \log \pi_{\theta}(a | s) Q(s, a)]
\end{aligned}$$

C) Let $\nabla_{\theta} J(\theta) = \mathbb{E}_{P_{\pi}, \pi} [g]$ where

$$g := \nabla_{\theta} \log \pi_{\theta}(a | s) (Q(s, a) - b(s))$$

$$\begin{aligned}
\text{Var}(g) &= E_{\pi, \pi}[(g - E_{\pi, \pi}[g])^T(g - E_{\pi, \pi}[g])] \\
&= E_{\pi, \pi}[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s)] b(s)^2 \\
&\quad - 2 E_{\pi, \pi}[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s) \hat{Q}(s, a)] b(s) \\
&\quad + E_{\pi, \pi}[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s) \hat{Q}(s, a)^2] \\
&\quad - E_{\pi, \pi}[\nabla_{\theta} \log \pi(a|s)(\hat{Q}(s, a) - b(s))]^T \\
&\quad E_{\pi, \pi}[\nabla_{\theta} \log \pi(a|s)(\hat{Q}(s, a) - b(s))]
\end{aligned}$$

from one hand:

$$\begin{aligned}
E_{\pi, \pi}[\nabla_{\theta} \log \pi(a|s) b(s)] &= E_{\pi} [E_{\pi | \pi}[\nabla_{\theta} \log \pi(a|s) b(s) | s]] \\
&= E_{\pi} [\int_a \nabla_{\theta} \log \pi(a|s) b(s) da | s] \\
&= E_{\pi} [b(s) \nabla_{\theta} \int_a \log \pi(a|s) da | s] \\
&= E_{\pi} [b(s) \nabla_{\theta} | | s] \\
&= E_{\pi} [0 | s] \\
&= 0
\end{aligned}$$

Therefore:

$$\begin{aligned}
\text{Var}(g) &= E_{\pi, \pi}[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s)] b(s)^2 \\
&\quad - 2 E_{\pi, \pi}[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s) \hat{Q}(s, a)] b(s)
\end{aligned}$$

$$\Rightarrow \frac{\partial}{\partial b(s)} \text{Var}(g) = 2 E_{\pi}[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s)] b^*(s) \\
- 2 E_{\pi}[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s) \hat{Q}(s, a)] = 0$$

$$\implies b^*(s) = \frac{E_n[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s) \hat{Q}(s, a)]}{E_n[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s)]}$$

D) Let's assume $\nabla_{\theta} \log \pi(a|s)$ and $\hat{Q}(s, a)$ are independent, therefore:

$$\begin{aligned} b^*(s) &= \frac{E_n[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s) \hat{Q}(s, a)]}{E_n[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s)]} \\ &= \frac{E_n[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s)] E_{\pi(s)}[\hat{Q}(s, a)]}{E_n[\nabla_{\theta} \log \pi(a|s)^T \nabla_{\theta} \log \pi(a|s)]} \\ &= E_{\pi}[\hat{Q}(s, a)] \\ &= V^{\pi}(s) \end{aligned}$$

(Q3)

A) Theorem: Let μ and ν be two probability distribution functions on \mathcal{X} . Then if $\forall x \in \mathcal{X}, |\mu(x) - \nu(x)| = \epsilon$, there is a coupling $P(X, Y)$, such that $P(X \neq Y) = \epsilon$.

Proof: Let

$$\begin{aligned} p &= \sum_{x \in \mathcal{X}} \mu(x) \wedge \nu(x) = \sum_{\substack{x \in \mathcal{X} \\ \mu(x) \leq \nu(x)}} \mu(x) + \sum_{\substack{x \in \mathcal{X} \\ \mu(x) > \nu(x)}} \nu(x) \\ &= \sum_{\substack{x \in \mathcal{X} \\ \mu(x) \leq \nu(x)}} \mu(x) + \sum_{\substack{x \in \mathcal{X} \\ \mu(x) > \nu(x)}} \nu(x) + \sum_{\substack{x \in \mathcal{X} \\ \mu(x) > \nu(x)}} \mu(x) - \sum_{\substack{x \in \mathcal{X} \\ \mu(x) > \nu(x)}} \mu(x) \\ &= 1 - \left(\sum_{\substack{x \in \mathcal{X} \\ \mu(x) > \nu(x)}} \mu(x) - \nu(x) \right) \\ &= 1 - |\mu(X) - \nu(X)| \end{aligned}$$

because let $B = \{x : \mu(x) \geq \nu(x)\}$ and let $A \subset \mathcal{X}$ be any event. then

$$\mu(A) - \nu(A) \leq \mu(A \cap B) - \nu(A \cap B) \leq \mu(B) - \nu(B)$$

Since for any $x \in A \cap B^c$, $\mu(x) - \nu(x)$, so difference of probability without these terms is greater, and hence the first inequality. By the same argument, the second inequality can be prove. and in the same way we have

$$\nu(A) - \mu(A) \leq \nu(B^c) - \mu(B^c)$$

therefore $\mu(B) - \nu(B) = \nu(B^c) - \mu(B^c)$.

Now consider a coin with the probability of heads equal to p . Flip the coin

(1) If it comes up heads, then choose a random value z by the probability distribution:

$$\gamma_1(z) = \frac{\mu(z) \wedge \nu(z)}{p}$$

and set $X = Y = z$.

(2) else, choose x according the probability distribution

$$\gamma_2(x) = \begin{cases} \frac{\mu(x) - \nu(x)}{1-p} & , \text{ if } \mu(x) > \nu(x) \\ 0 & , \text{ otherwise} \end{cases}$$

and Set $X = x$. In a same way choose y according the probability distribution

$$\gamma_3(y) = \begin{cases} \frac{\nu(y) - \mu(y)}{1-p} & , \text{ if } \nu(y) > \mu(y) \\ 0 & , \text{ otherwise} \end{cases}$$

and set $Y = y$. Obviously in this coupling, we have:

$$P\{X = Y\} = p = 1 - \varepsilon, \quad P\{X \neq Y\} = 1 - p = \varepsilon$$

and

$$p\gamma_1 + (1-p)\gamma_2 = \mu$$

$$p\gamma_1 + (1-p)\gamma_3 = \nu$$

thus the marginals of $P(X, Y)$, match μ and ν .

Now consider the above coupling for distributions $\pi_{\theta}(a_t | s_t)$ and $\pi_{\theta'}(a_t | s_t)$, where $|\pi_{\theta}(a_t | s_t) - \pi_{\theta'}(a_t | s_t)| \leq \varepsilon$. By this coupling, the probability of taking different actions in each state s_t by π_{θ} and $\pi_{\theta'}$ would be at most ε . Thus for a state s_t ,

$$P_{\theta'}(s_t) = P_{nm}^t P_{\theta}(s_t) + (1 - P_{nm}^t) P_m(s_t)$$

$$P_{nm} = 1 - \varepsilon \quad (\text{probability of making no mistakes})$$

P_m = probability of some distribution other than P_{θ} .

So

$$|P_{\theta'}(s_t) - P_{\theta}(s_t)| = (1 - (1 - \varepsilon)^t) |P_m(s_t) - P_{\theta}(s_t)|$$

By the Bernoulli identity, $(1 - \varepsilon)^t \geq 1 - \varepsilon t$. Thus

$$|P_{\theta'}(s_t) - P_{\theta}(s_t)| \leq \varepsilon t |P_m(s_t) - P_{\theta}(s_t)|$$

Also for any event $A \in \mathcal{X}$, we have

$$\begin{aligned} P_m(A) - P_{\theta}(A) &= P\{X \in A\} - P\{Y \in A\} \\ &\leq P\{X \in A, Y \notin A\} \\ &\leq P(X \neq Y) \\ &\leq 1 \end{aligned}$$

$$\Rightarrow |P_m(s_t) - P_{\theta}(s_t)| = 2 |P_m(s_t) - P_{\theta}(s_t)|_{TV} \leq 2$$

and

$$|P_{\theta'}(s_t) - P_\theta(s_t)| \leq 2\varepsilon t$$

B) Now assume $|\pi_{\theta'}(a_t | s_t) - \pi_\theta(a_t | s_t)| \leq \varepsilon$ for all s_t . therefore $|P_{\theta'}(s_t) - P_\theta(s_t)| \leq 2\varepsilon t$ for all s_t . Now consider an arbitrary function of s_t , such as $f(s_t)$. We have:

$$\begin{aligned} E_{P_{\theta'}}[f(s_t)] &= \sum_{s_t} P_{\theta'}(s_t) f(s_t) \\ &= \sum_{s_t} (P_{\theta'}(s_t) + P_\theta(s_t) - P_\theta(s_t)) f(s_t) \\ &= \sum_{s_t} P_\theta(s_t) f(s_t) - \sum_{s_t} (P_\theta(s_t) - P_{\theta'}(s_t)) f(s_t) \\ &\geq \sum_{s_t} P_\theta(s_t) f(s_t) - |P_\theta(s_t) - P_{\theta'}(s_t)| \max_{s_t} f(s_t) \\ &\geq E_{P_\theta}[f(s_t)] - 2\varepsilon t \max_{s_t} f(s_t) \end{aligned}$$

By setting $f(s_t) = E_{a_t \sim \pi_\theta(a_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_\theta(a_t | s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t) \right]$ for infinite horizon we obtain:

$$\begin{aligned} \sum_{t=0}^T E_{P_{\theta'}(s_t)} \left[E_{a_t \sim \pi_\theta(a_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_\theta(a_t | s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t) \right] \right] \\ \geq \sum_{t=0}^T E_{P_\theta(s_t)} \left[E_{a_t \sim \pi_\theta(a_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_\theta(a_t | s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t) \right] \right] - \sum_{t=0}^T 2\varepsilon t C \\ \geq \sum_{t=0}^T E_{P_\theta(s_t)} \left[E_{a_t \sim \pi_\theta(a_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_\theta(a_t | s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t) \right] \right] - \sum_{t=0}^T 2\varepsilon t C \end{aligned}$$

which $\sum_{t=0}^T 2\varepsilon t c = 2x \frac{T(T+1)}{2} \varepsilon c$. When $c \leq r_{\max}$ so

$$\begin{aligned} & \sum_{t=0}^T \underbrace{\mathbb{E}_{\theta^{(s_t)}}}_{P_\theta(s_t)} \left[\mathbb{E}_{a_t \sim \pi_\theta(s_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_\theta(a_t | s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t) \right] \right] \\ & \geq \sum_{t=0}^T \mathbb{E}_{\theta^{(s_t)}} \left[\mathbb{E}_{a_t \sim \pi_\theta(s_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_\theta(a_t | s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t) \right] \right] - O(T^2 r_{\max}) \end{aligned}$$

And for infinite horizon:

$$\begin{aligned} & \sum_t \mathbb{E}_{\theta^{(s_t)}} \left[\mathbb{E}_{a_t \sim \pi_\theta(s_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_\theta(a_t | s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t) \right] \right] \\ & \geq \sum_t \mathbb{E}_{\theta^{(s_t)}} \left[\mathbb{E}_{a_t \sim \pi_\theta(s_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_\theta(a_t | s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t) \right] \right] - \sum_t 2\varepsilon t c \end{aligned}$$

$$c(t) = \max_{s_t, a_t} |\gamma^t A^{\pi_\theta}(s_t, a_t)| \leq \gamma^t \max_{s, a} |A^{\pi_\theta}(s, a)| = \gamma^t c$$

$$\begin{aligned} \sum_t 2\varepsilon t c(t) & \leq \sum_t 2\varepsilon t \gamma^t c \approx \frac{d}{dt} \sum_t 2\varepsilon \gamma^{t+1} c \\ & = 2\varepsilon c \frac{d}{dt} \left(\frac{1}{1-\gamma} - 1 \right) \\ & = \frac{2\varepsilon c}{(1-\gamma)^2} = O\left(\frac{r_{\max}}{(1-\gamma)^2}\right) \end{aligned}$$

$$\begin{aligned} & \sum_t \mathbb{E}_{\theta^{(s_t)}} \left[\mathbb{E}_{a_t \sim \pi_\theta(s_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_\theta(a_t | s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t) \right] \right] \\ & \geq \sum_t \mathbb{E}_{\theta^{(s_t)}} \left[\mathbb{E}_{a_t \sim \pi_\theta(s_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_\theta(a_t | s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t) \right] \right] - O\left(\frac{r_{\max}}{(1-\gamma)^2}\right) \end{aligned}$$

C) By the Pinsker's inequality:

$$|\mathcal{H}_{\theta'}(a_t|s_t) - \mathcal{H}_{\theta}(a_t|s_t)|_{TV} \leq \sqrt{\frac{1}{2} D_{KL}(\pi_{\theta}(a_t|s_t) || \pi_{\theta'}(a_t|s_t))}$$

$$\Rightarrow 2|\mathcal{H}_{\theta'}(a_t|s_t) - \mathcal{H}_{\theta}(a_t|s_t)|_{TV}^2 \leq D_{KL}(\pi_{\theta}(a_t|s_t) || \pi_{\theta'}(a_t|s_t))$$

D) Let $\sum_t E_{P_{\theta}}[E_{\pi}[\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta'}(a_t|s_t)} \gamma^t A^{\pi_{\theta}}(s_t, a_t)]] = \bar{A}(\theta')$. So the

Taylor's approximation of $\bar{A}(\theta')$ is:

$$\bar{A}(\theta') \approx \bar{A}(\theta) + \nabla_{\theta} \bar{A}(\theta)^T (\theta' - \theta)$$

$$\begin{aligned} \nabla_{\theta} \bar{A}(\theta) &= \nabla_{\theta} \sum_t E_{P_{\theta}}[E_{\pi}[\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta'}(a_t|s_t)} \gamma^t A^{\pi_{\theta}}(s_t, a_t)]] \\ &= \sum_t E_{P_{\theta}}[E_{\pi}[\gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) A^{\pi_{\theta}}(s_t, a_t)]] \\ &= \nabla_{\theta} J(\theta) \end{aligned}$$

Therefore:

$$\bar{A}(\theta') \approx \sum_t E_{P_{\theta}}[E_{\pi}[\gamma^t A^{\pi_{\theta}}(s_t, a_t)]] + \nabla_{\theta} J(\theta)^T (\theta' - \theta)$$

Q4

We know $Q^{\mu_\theta}(s, \mu_\theta(s)) = E_{a \sim \mu_\theta(s)}[Q^{\mu_\theta}(s, a)]$. Thus

$$\begin{aligned}
\nabla_\theta V^{\mu_\theta}(s) &= \nabla_\theta Q^{\mu_\theta}(s, \mu_\theta(s)) \\
&= \nabla_\theta (\gamma r(s, \mu_\theta(s)) + \gamma E_{s' \sim p(s'|s, \mu_\theta(s))}[V^{\mu_\theta}(s')]) \\
&= \nabla_\theta (\gamma r(s, \mu_\theta(s)) + \int_S \gamma p(s'|s, \mu_\theta(s)) V^{\mu_\theta}(s') ds') \\
&= \nabla_\theta \mu_\theta(s) \nabla_a r(s, a)|_{a=\mu_\theta(s)} + \int_S \gamma p(s'|s, \mu_\theta(s)) \nabla_\theta V^{\mu_\theta}(s') ds' \\
&\quad + \int_S \nabla_\theta \mu_\theta(s) \nabla_a p(s'|s, a)|_{a=\mu_\theta(s)} V^{\mu_\theta}(s') ds' \quad (\text{chain rule}) \\
&= \nabla_\theta \mu_\theta(s) \nabla_a (r(s, a) + \int_S \gamma p(s'|s, a) V^{\mu_\theta}(s') ds')|_{a=\mu_\theta(s)} \\
&\quad + \int_S \gamma p(s'|s, \mu_\theta(s)) \nabla_\theta V^{\mu_\theta}(s') ds' \\
&= \nabla_\theta \mu_\theta(s) \nabla_a Q^{\mu_\theta}(s, a)|_{a=\mu_\theta(s)} + \int_S \gamma p(s'|s, \mu_\theta(s)) \nabla_\theta V^{\mu_\theta}(s') ds' \\
&= \nabla_\theta \mu_\theta(s) \nabla_a Q^{\mu_\theta}(s, a)|_{a=\mu_\theta(s)} + \int_S \gamma p(s \rightarrow s', 1, \mu_\theta) \nabla_\theta V^{\mu_\theta}(s') ds' \\
&= \nabla_\theta \mu_\theta(s) \nabla_a Q^{\mu_\theta}(s, a)|_{a=\mu_\theta(s)} + \int_S \gamma p(s \rightarrow s', 1, \mu_\theta) \nabla_\theta V^{\mu_\theta}(s') ds'
\end{aligned}$$

the s is compact
 and $\|\nabla_\theta V^{\mu_\theta}(s)\|$
 and all the other
 ∇s are bound

By recursively applying the above relation we will have:

$$\begin{aligned}
\nabla_\theta V^{\mu_\theta}(s) &= \nabla_\theta \mu_\theta(s) \nabla_a Q^{\mu_\theta}(s, a)|_{a=\mu_\theta(s)} + \int_S \gamma p(s \rightarrow s', 1, \mu_\theta) \nabla_\theta V^{\mu_\theta}(s') ds' \\
&= \nabla_\theta \mu_\theta(s) \nabla_a Q^{\mu_\theta}(s, a)|_{a=\mu_\theta(s)} \\
&\quad + \int_S \gamma p(s \rightarrow s', 1, \mu_\theta) \nabla_\theta \mu_\theta(s') \nabla_a Q^{\mu_\theta}(s', a)|_{a=\mu_\theta(s')} ds' \\
&\quad + \int_S \gamma p(s \rightarrow s', 1, \mu_\theta) \int_S \gamma p(s' \rightarrow s'', 1, \mu_\theta) \nabla_\theta V^{\mu_\theta}(s'') ds'' ds'
\end{aligned}$$

$$\begin{aligned}
&= \nabla_{\theta} \mu_{\theta}(s) \nabla_a Q^{\theta}(s, a) \Big|_{a=\mu_{\theta}(s)} \\
&\quad + \int_S \gamma p(s \rightarrow s', 1, \mu_{\theta}) \nabla_{\theta} \mu_{\theta}(s') \nabla_a Q^{\mu_{\theta}}(s', a) \Big|_{a=\mu_{\theta}(s')} ds' \\
&\quad + \int_S \gamma^2 p(s \rightarrow s', 2, \mu_{\theta}) \nabla_{\theta} V^{\mu_{\theta}}(s') ds' \\
&\vdots \\
&= \int_S \sum_{t=0}^{\infty} \gamma^t p(s \rightarrow s', t, \mu_{\theta}) \nabla_{\theta} \mu_{\theta}(s') \nabla_a Q(s', a) \Big|_{a=\mu_{\theta}(s')} ds'
\end{aligned}$$

Wow we have

$$\begin{aligned}
\nabla_{\theta} J(\mu_{\theta}) &= \nabla_{\theta} E \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) / \mu \right] \\
&= \nabla_{\theta} E [V^{\mu_{\theta}}(s)] \\
&= \nabla_{\theta} \int_S p_{\mu}(s) V^{\mu_{\theta}}(s) ds \\
&= \int_S p_{\mu}(s) \nabla_{\theta} V^{\mu_{\theta}}(s) ds \\
&= \int_S \int_S \sum_{t=0}^{\infty} \gamma^t p_{\mu}(s) p(s \rightarrow s', t, \mu_{\theta}) \nabla_{\theta} \mu_{\theta}(s') \nabla_a Q(s', a) \Big|_{a=\mu_{\theta}(s')} ds' ds \\
&= \int_S p_{\mu}(s) \nabla_{\theta} \mu_{\theta}(s) \nabla_a Q(s, a) \Big|_{a=\mu_{\theta}(s)} ds \\
&= E_{s \sim p_{\mu}} [\nabla_{\theta} \mu_{\theta}(s) \nabla_a Q(s, a) \Big|_{a=\mu_{\theta}(s)}]
\end{aligned}$$