

(Alireza Gargoori Motlagh - 98102176)

((Homework 2 - Reinforcement Learning))

1 Policy Gradient

a) As we will prove in part (c) (also proved in class!), we have:

$$J(\theta) = \mathbb{E}_{\tau \sim P_\theta(\tau)} \left[\sum_t \gamma^t r(s_t, a_t) \right] = \mathbb{E}_{\tau \sim P_\theta(\tau)} \left[\gamma^t r(\tau) \right]$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim P_\theta(\tau)} \left[\nabla_{\theta} \log P_\theta(\tau) r(\tau) \right]$$

where

$$P_\theta(\tau_i) = P(s_{i,0}) \prod_{t=1}^T \pi(a_{i,t} | s_{i,t}) P(s_{i,t+1} | a_{i,t}, s_{i,t})$$

$$r(\tau_i) = \sum_{t'=t}^T \gamma^{t'} r(s_{i,t'}, a_{i,t'})$$

$$\Rightarrow \nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim P_\theta(\tau)} \left[\nabla_{\theta} \log \pi_{\theta}(\tau) (r(\tau) - b(s)) \right] = \dots$$

$$= \mathbb{E}_{\tau \sim P_\theta(\tau)} \left[\nabla_\theta \log P_\theta(\tau) r(\tau) \right] - \mathbb{E}_{\tau \sim P_\theta(\tau)} \left[\nabla_\theta \log P_\theta(\tau) b(s) \right]$$

$$* \mathbb{E}_{\tau \sim P_\theta(\tau)} \left[\nabla_\theta \log P_\theta(\tau) b(s) \right] = \int P_\theta(\tau) \nabla_\theta \log P_\theta(\tau) b(s) d\tau =$$

$$\int P_\theta(\tau) \frac{\nabla_\theta P_\theta(\tau)}{P_\theta(\tau)} b(s) d\tau = \int \nabla_\theta P_\theta(\tau) b(s) d\tau = b(s) \int \nabla_\theta P_\theta(\tau) d\tau =$$

$$b(s) \nabla \int P_\theta(\tau) d\tau = b \nabla_\theta 1 = 0 \Rightarrow$$

$$\nabla_\theta J_b(\theta) = \mathbb{E}_{\tau \sim P_\theta(\tau)} \left[\nabla_\theta \log P_\theta(\tau) r(\tau) \right] = \nabla_\theta J(\theta)$$

The above formula shows that $J(\theta)$ and the baseline - version $J_b(\theta)$

when using the true distribution $P(\tau)$ are the same ; therefore their

expectation by averaging would be the same and $\nabla_\theta J_b(\theta)$ is unbiased.

Therefore , subtracting a baseline (which only depends on states) is

unbiased in expectation .

$$(b) \quad \nabla_{\theta} J_b(\theta) = \mathbb{E}_{\tau \sim P_{\theta}(\tau)} \left[\nabla_{\theta} \log p_{\theta}(\tau) (r(\tau) - b(s)) \right]$$

$$\text{Var}(\nabla_{\theta} J_b(\theta)) = \mathbb{E}_{\tau \sim P_{\theta}(\tau)} \left[(\nabla_{\theta} \log p_{\theta}(\tau) (r(\tau) - b))^2 \right] - \left(\mathbb{E}_{\tau \sim P_{\theta}(\tau)} (\nabla_{\theta} \log p_{\theta}(\tau) (r(\tau) - b)) \right)^2$$

$$\stackrel{(a)}{=} \left(\mathbb{E}_{\tau \sim P_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)] \right)^2$$

(not a function of b)

$$\Rightarrow \frac{\partial \text{Var}(\nabla_{\theta} J_b(\theta))}{\partial b} = \frac{\partial}{\partial b} \mathbb{E}_{\tau \sim P_{\theta}(\tau)} \left[(\nabla_{\theta} \log p_{\theta}(\tau) (r(\tau) - b))^2 \right] =$$

$$\cancel{\frac{\partial}{\partial b} \left(\mathbb{E}_{\tau} \left[(\nabla_{\theta} \log p_{\theta})^2 r(\tau)^2 \right] - 2 \mathbb{E}_{\tau} \left[(\nabla_{\theta} \log p_{\theta})^2 r(\tau) b \right] + \mathbb{E}_{\tau} \left[(\nabla_{\theta} \log p_{\theta})^2 b^2 \right] \right)}$$

$$= -2 \mathbb{E}_{\tau} \left[(\nabla_{\theta} \log p_{\theta}(\tau))^2 r(\tau) \right] + 2b \mathbb{E}_{\tau} \left[(\nabla_{\theta} \log p_{\theta}(\tau))^2 \right] = 0 \Rightarrow$$

$$b^* = \frac{\mathbb{E}_{\tau \sim P_{\theta}(\tau)} \left[(\nabla_{\theta} \log p_{\theta}(\tau))^2 r(\tau) \right]}{\mathbb{E}_{\tau \sim P_{\theta}(\tau)} \left[(\nabla_{\theta} \log p_{\theta}(\tau))^2 \right]}$$

which is the weighted version of expected rewards by the magnitude

of their log probability. (note that: $\nabla_{\theta} \log p_{\theta}(\tau) = \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$)

$$\begin{aligned}
 (c) \quad V_{(s)}^{\pi_\theta} &= \sum_a \pi_\theta(a|s) Q_{\pi_\theta}^{\pi_\theta}(s, a) \Rightarrow \\
 \nabla_\theta V_{(s)}^{\pi_\theta} &= \sum_a \nabla_\theta \pi_\theta(a|s) Q_{\pi_\theta}^{\pi_\theta}(s, a) + \sum_a \pi_\theta(a|s) \nabla_\theta Q_{\pi_\theta}^{\pi_\theta}(a, s) \\
 &= \sum_a \nabla_\theta \pi_\theta(a|s) Q_{\pi_\theta}^{\pi_\theta}(s, a) + \sum_a \pi_\theta(a|s) \nabla_\theta \sum_{s'} P(s'|s, a) (r + V_{(s')}^{\pi_\theta}) \\
 &= \sum_a \nabla_\theta \pi_\theta(a|s) Q_{\pi_\theta}^{\pi_\theta}(s, a) + \sum_a \pi_\theta(a|s) \sum_{s'} P(s'|s, a) \nabla_\theta V_{(s')}^{\pi_\theta} \\
 \text{iterate} \quad &= \sum_a \nabla_\theta \pi_\theta(a|s) Q_{\pi_\theta}^{\pi_\theta}(s, a) + \sum_a \pi_\theta(a|s) \sum_{s'} P(s'|s, a) \times \\
 &\quad \left(\sum_{a'} \nabla_\theta \pi_\theta(a'|s') Q_{\pi_\theta}^{\pi_\theta}(s', a') + \sum_{a'} \pi_\theta(a'|s') \sum_{s''} P(s''|s', a') \nabla_\theta V_{(s'')}^{\pi_\theta} \right) \\
 \text{iteration} \quad &= \dots = \\
 &= \sum_{s'} \sum_{t=0}^{\infty} P(s \rightarrow s', t, \pi_\theta) \sum_a \nabla_\theta \pi_\theta(a|s') Q_{\pi_\theta}^{\pi_\theta}(s', a)
 \end{aligned}$$

$\Pr(s \rightarrow s', t, \pi)$ is the probability of transitioning from s to s'
 with t steps under policy π_θ . (The same as Problem 3a and
 SB-RL Book). Also, using: $\nabla_\theta \pi_\theta(a|s) = \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s)$

would be discounted

$$\begin{aligned}
 \nabla_\theta V_{(s)}^{\pi_\theta} &= \sum_s \sum_{t=0}^{\infty} \underbrace{\Pr(s \rightarrow s, t, \pi_\theta)}_{\text{discounted}} \sum_a \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) Q_{\pi_\theta}^{\pi_\theta}(s, a) \\
 &= \sum_s \sum_{t=0} \gamma^t \Pr(s_t = s | s, \pi_\theta) \sum_a \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) Q_{\pi_\theta}^{\pi_\theta}(s, a)
 \end{aligned}$$

Taking expectation on initial states distribution:

$$\Rightarrow \nabla_{\theta} V(\mu) = \nabla_{\theta} \mathbb{E}_{s_0 \sim \mu} [V(s_0)] = \mathbb{E}_{s_0 \sim \mu} [\nabla_{\theta} V^{\pi_{\theta}}(s_0)] =$$

$$\sum_{s_0} \mu(s_0) \sum_s \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s | s_0, \pi_{\theta}) \sum_a \pi_{\theta}(a_t | s_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q^{\pi_{\theta}}(s_t, a_t)$$

$$= \mathbb{E}_{\mu, \pi} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q^{\pi_{\theta}}(a_t | s_t) \right] \xrightarrow{\tau \sim \Pr_{\mu}^{\pi}}$$

$$\boxed{\nabla_{\theta} V^{\pi_{\theta}}(\mu) = \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q^{\pi_{\theta}}(a_t | s_t) \right]}$$

$$(d) \quad \mathbb{E}_{\tau \sim \Pr_{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right] = \frac{1}{1-\gamma} \mathbb{E}_{s_0 \sim d^{\pi_{\theta}}_{s_0}} \mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)} [f(s, a)]$$

Substituting $f(s_t, a_t)$ with: $\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q^{\pi_{\theta}}(a_t | s_t)$ in the

Above equation yields:

$$\nabla_{\theta} V^{\pi}(\mu) = \mathbb{E}_{\pi \sim P_{\mu}^{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi(a_t | s_t) Q^{\pi_{\theta}}(s_t, a_t) \right] =$$

$$\frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta}}} \mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)} \left[\nabla_{\theta} \log \pi_{\theta}(a | s) Q^{\pi_{\theta}}(s, a) \right] \implies$$

$$\boxed{\nabla_{\theta} V^{\pi_{\theta}}(\mu) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta}}} \mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)} \left[\nabla_{\theta} \log \pi_{\theta}(a | s) Q^{\pi_{\theta}}(s, a) \right]}$$

As we proved in part (a), $\nabla_{\theta} J(\theta) = \nabla_{\theta} V^{\pi}(\mu)$ is unbiased

in expectation when subtracting a baseline which only

depends on states, e.g. Value of state $V^b(s)$! ; therefore:

$$\nabla_{\theta} V^{\pi}(\mu) = \mathbb{E}_{\pi \sim P_{\mu}^{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi(a_t | s_t) Q^{\pi_{\theta}}(s_t, a_t) \right] =$$

$$\nabla_{\theta} V^{\pi}(\mu) = \mathbb{E}_{\pi \sim P_{\mu}^{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi(a_t | s_t) \left[Q^{\pi_{\theta}}(s_t, a_t) - V^{\pi_{\theta}}(s_t) \right] \right] =$$

$$\mathbb{E}_{\pi \sim P_{\mu}^{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi(a_t | s_t) A^{\pi_{\theta}}(s_t, a_t) \right] \xrightarrow[\text{(8)}]{\text{Using eq.}}$$

$$\boxed{\nabla_{\theta} V^{\pi_{\theta}}(\mu) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta}}} \mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)} \left[\nabla_{\theta} \log \pi_{\theta}(a | s) A^{\pi_{\theta}}(s, a) \right]}$$

2 Value-based Methods for Continuous Actions

- a) Since Q-learning is based on a greedy policy which involves finding the $\left[\arg \max_a Q(s, a) \right] \neq \max_a Q(s, a)$ and it requires an iterative optimization process at every step. In other words, optimizing in the inner loop of an optimization problem would lead to massive increase in the computation time.
- Also, discretizing the action space is limited as well and the number of actions increases exponentially with the degrees of freedom. (DDPG original paper)

(b)

1. $Q_\phi(s, a) = -\frac{1}{2} (a - \mu_\theta(s))^T P_\phi(s) (a - \mu_\theta(s)) + V_\phi(s)$

$$\arg \max_a Q_\phi(s, a) : \frac{\partial Q_\phi}{\partial a} = 0 \Rightarrow -\frac{1}{2} \times 2 \times P_\phi(s) (a - \mu_\theta(s)) = 0$$

$$\Rightarrow P_{\emptyset}^{(s)} (a - \mu_{\emptyset}^{(s)}) = 0 \xrightarrow{*} a = \mu_{\emptyset}^{(s)}$$

* It can have other answers if $P_{\emptyset}^{(s)}$ is not positive definite.

$$\Rightarrow Q_{\emptyset}^{(s,a)} \Big|_{a=\mu_{\emptyset}^{(s)}} = 0 + V_{\emptyset}^{(s)} = V_{\emptyset}^{(s)}$$

$$\begin{aligned} &\Rightarrow \boxed{\arg \max_a Q_{\emptyset}^{(s,a)} = \mu_{\emptyset}^{(s)}} \\ &\quad \boxed{\max_a Q_{\emptyset}^{(s,a)} = V_{\emptyset}^{(s)}} \end{aligned}$$

Using simple functions for state-action values, would result in a model which cannot capture the complex behaviour needed to execute in the environment ; however , it results in a simpler model which needs less computation & etc.
 (Somehow the same as underfitting !)

2.a DDPG is a model-free, off-policy, actor-critic method

which alleviates the problem of a Q-learning in

Continuous action spaces by interleaving learning an approximator

to $Q^*(s, a)$ with learning an approximator to $a^*(s) = \mu_\theta(s)$.

Since the action space is Continuous, the function $Q^*(s, a)$

is presumed to be differentiable and we can approximate

$$\max_a Q(s, a) \approx Q(s, \mu_\theta(s)) .$$

The Q-learning part minimizes (Critic):

$$L(\phi, D) = \mathbb{E}_{M \sim D} \left[(Q_\phi(s, a) - (r + \gamma \max_{a' \in \Phi_{\text{targ}}} Q_\phi(s', \mu_{\theta_{\text{targ}}}(s'))))^2 \right]$$

where D is an experience replay buffer (we are able to do this

sense Bellman eq. doesn't care which transition tuples are

used and etc, Since $Q^*(s, a)$ should satisfy Bellman optimality

for all transitions.) Also, target networks are used for

Stability of learning by updating the target network by polyak

Averaging once per main network update: ($\tau \ll 1$)

$$\phi_{\text{targ}} \leftarrow \tau \phi + (1 - \tau) \phi_{\text{targ}}$$

$$\theta_{\text{targ}} \leftarrow \tau \theta + (1 - \tau) \theta_{\text{targ}}$$

The Objective function for actor is also:

$$\max_{\theta} \mathbb{E}_{\substack{s \sim D \\ \phi}} [Q(s, \mu_{\theta}(s))]$$

Algorithm 1 DDPG algorithm

Randomly initialize critic network $Q(s, a | \theta^Q)$ and actor $\mu(s | \theta^\mu)$ with weights θ^Q and θ^μ .

Initialize target network Q' and μ' with weights $\theta^{Q'} \leftarrow \theta^Q, \theta^{\mu'} \leftarrow \theta^\mu$

Initialize replay buffer R

for episode = 1, M **do**

 Initialize a random process \mathcal{N} for action exploration

 Receive initial observation state s_1

for t = 1, T **do**

 Select action $a_t = \mu(s_t | \theta^\mu) + \mathcal{N}_t$ according to the current policy and exploration noise

 Execute action a_t and observe reward r_t and observe new state s_{t+1}

 Store transition (s_t, a_t, r_t, s_{t+1}) in R

 Sample a random minibatch of N transitions (s_i, a_i, r_i, s_{i+1}) from R

 Set $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1} | \theta^{\mu'}) | \theta^{Q'})$

 Update critic by minimizing the loss: $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i | \theta^Q))^2$

 Update the actor policy using the sampled gradient:

$$\nabla_{\theta^\mu} \mu|_{s_i} \approx \frac{1}{N} \sum_i \nabla_a Q(s, a | \theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s | \theta^\mu)|_{s_i}$$

 Update the target networks:

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$$

$$\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$$

end for

end for

2.b In the DDPG algorithm, the actor and critic work with each other in the sense that the evaluation of critic is passed to the actor to maximize the performance through policy gradient and therefore, passing of gradients from critic to actor provides a better evaluation of value of the action.

However, in the classic REINFORCE algorithm, the trajectories are obtained from the current policy and the estimation of value of state has a high variance and values are not learnt through any critic network which results in a poor estimate and subsequent problems. Therefore actor-critic methods such as DDPG alleviate the problem of high variance estimate of value of states.

$$2.C \quad Q^{\mu}(s, \mu_{\theta}(s)) = r(s, \mu_{\theta}(s)) + \int_{\mathcal{S}} \gamma P(s' | s, \mu_{\theta}(s)) V^{\mu}(s') ds'$$

$$\Rightarrow \nabla_{\theta} V^{\mu}(s) = \nabla_{\theta} Q^{\mu}(s, \mu_{\theta}(s)) = \left. \nabla_{\theta} \mu_{\theta}(s) \nabla_a r(s, \mu_{\theta}(s)) \right|_{a=\mu_{\theta}(s)} +$$

Chain rule

$$\int_{\mathcal{S}} \gamma \nabla_{\theta} \mu_{\theta}(s) \nabla_a P(s' | s, a) \Big|_{a=\mu_{\theta}(s)} V^{\mu}(s') ds' + \int_{\mathcal{S}} \gamma P(s' | s, \mu_{\theta}(s)) \nabla_{\theta} V^{\mu}(s') ds'$$

$$= \nabla_{\theta} \mu_{\theta}(s) \nabla_a \left(r(s, a) + \int_{\mathcal{S}} \gamma P(s' | s, a) V^{\mu}(s') ds' \right) \Big|_{a=\mu_{\theta}(s)} +$$

$\underbrace{Q^{\mu}(s, a)}$

$$\int_{\mathcal{S}} \gamma P(s' | s, \mu_{\theta}(s)) \nabla_{\theta} V^{\mu}(s') ds' = \nabla_{\theta} \mu_{\theta}(s) \nabla_a Q^{\mu}(s, a) \Big|_{a=\mu_{\theta}(s)} +$$

$$\int_{\mathcal{S}} \gamma P(s' | s, \mu_{\theta}(s)) \nabla_{\theta} V^{\mu}(s') ds' = \nabla_{\theta} \mu_{\theta}(s) \nabla_a Q^{\mu}(s, a) \Big|_{a=\mu_{\theta}(s)} +$$

$\underbrace{\text{Iterate}}$

$$\int_{\mathcal{S}} \gamma \Pr(s \rightarrow s', 1, \mu_{\theta}) \nabla_{\theta} \mu_{\theta}(s) \nabla_a Q^{\mu}(s, a) \Big|_{a=\mu_{\theta}(s')} +$$

$$\int_{\mathcal{S}} \gamma \Pr(s \rightarrow s', 1, \mu_{\theta}) \int_{\mathcal{S}} \gamma \Pr(s' \rightarrow s'', 1, \mu_{\theta}(s')) \nabla_{\theta} V^{\mu}(s'') ds'' ds' =$$

(Consider s' as the middle point to reach s'')

$$= \nabla_{\theta} \mu_{\theta}(s) \left| \nabla_a Q(s, a) \right|_{a=\mu_{\theta}(s)} + \int_s \gamma \Pr(s \rightarrow s', 1, \mu_{\theta}) \left| \nabla_{\theta} \mu_{\theta}(s') \nabla_a^{\mu} Q(s, a) \right|_{a=\mu_{\theta}(s')} ds'$$

$$+ \int_s \gamma^2 \Pr(s \rightarrow s'', 2, \mu_{\theta}) \underbrace{\nabla_{\theta} V^{\mu}(s'')}_{\text{iterate again}} ds'' = \dots$$

where $\Pr(s \rightarrow s', k, \mu_{\theta})$ is the probability of transition from state s to s' in k steps under policy μ_{θ} .

Since $\|\nabla_{\theta} V^{\mu}(s)\|$ is bounded, we can rewrite this :

$$\nabla_{\theta} V^{\mu}(s) = \int_s \sum_{t=0}^H \gamma^t \Pr(s \rightarrow s', t, \mu_{\theta}) \left| \nabla_{\theta} \mu_{\theta}(s') \nabla_a^{\mu} Q(s, a) \right|_{a=\mu_{\theta}(s')} ds'$$

Expectation over starting states distribution $p_o(s)$:

$$\begin{aligned} \nabla_{\theta} J(\mu_{\theta}) &= \nabla_{\theta} \mathbb{E}_{s \sim p_o} [V^{\mu}(s_o)] = \nabla_{\theta} \int_s p_o(s) V^{\mu}(s) ds = \int_s p_o(s) \nabla_{\theta} V^{\mu}(s) ds \\ &= \dots \end{aligned}$$

$$= \int_S p_\theta(s) \int_S \sum_{t=0}^H \gamma^t \Pr(s \rightarrow s', t, \mu_\theta) \nabla_\theta \mu_\theta(s') \nabla_a Q(s', a) |_{a=\mu_\theta(s')} ds' ds$$

$$= \int_S \int_S \sum_{t=0}^H \gamma^t p_\theta(s) \Pr(s \rightarrow s', t, \mu_\theta) \nabla_\theta \mu_\theta(s') \nabla_a Q(s', a) |_{a=\mu_\theta(s')} ds' ds$$

$$= \int_S \nabla_\theta \mu_\theta(s') \nabla_a Q(s', a) |_{a=\mu_\theta(s')} ds' \int_S \sum_{t=0}^H \gamma^t p_\theta(s) \Pr(s \rightarrow s', t, \mu_\theta) ds$$

$\underbrace{p^\mu(s)}$

$$= \int_S p^\mu(s) \nabla_\theta \mu_\theta(s) \nabla_a Q(s, a) |_{a=\mu_\theta(s)} ds'$$

where $p^\mu(s) = \sum_t \gamma^t p_\theta(s) \Pr(s \rightarrow s', t, \mu_\theta) ds$ is the discounted state distribution.

\Rightarrow

$$\nabla_\theta J(\mu_\theta) = \int_S p^\mu(s) \nabla_\theta \mu_\theta(s) \nabla_a Q(s, a) |_{a=\mu_\theta(s)}$$

$$= \mathbb{E}_{s \sim p^\mu} \left[\nabla_\theta \mu_\theta(s) \nabla_a Q(s, a) |_{a=\mu_\theta(s)} \right]$$

3 New Optimization method by Chaging Trust Region

$$(a) \quad V_{\tilde{\pi}}(s_0) = (1-\gamma) \mathbb{E}_{\rho, \tilde{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right] = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\rho, \tilde{\pi}} \left[(1-\gamma) R(s_t, a_t) \right]$$

$$= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\rho, \tilde{\pi}} \left[(1-\gamma) R(s_t, a_t) + V_{\pi}(s_t) - V_{\pi}(s_t) \right] =$$

$$\sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\rho, \tilde{\pi}} \left[(1-\gamma) R(s_t, a_t) + \gamma V_{\pi}(s_{t+1}) - V_{\pi}(s_t) \right] + V_{\pi}(s_0)$$

$$= V_{\pi}(s_0) + \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\tilde{\pi}, \rho} \left[A_{\pi}(s_t, a_t) \right] \Rightarrow$$

$$V_{\tilde{\pi}}(s_0) = V_{\pi}(s_0) + \frac{1}{1-\gamma} \mathbb{E}_{\tilde{\pi}, \rho} \left[A_{\pi}(s_t, a_t) \right] \xrightarrow{\mathbb{E}_{\rho}}$$

$\rho_{\tilde{\pi}}(s) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t=s | \tilde{\pi}, \rho)$

$$J(\tilde{\pi}) = J(\pi) + \iint_{S \times A} A_{\pi}^{\pi}(s, a) d\tilde{\pi}(a|s) d\rho_{\tilde{\pi}}(s)$$

(b) Original Optimization Problem:

$$\sup_{\tilde{\pi} \in \Pi} \left\{ \int_S \int_A A^\pi(s, a) d\tilde{\pi}(a|s) d\rho_\pi(s) \right\}$$

s.t. $\tilde{\pi} \in T_\epsilon := \left\{ \tilde{\pi} \in \Pi : \int_S C(\pi(.|s), \tilde{\pi}(.|s)) d\rho(s) \leq \epsilon \right\}$

We rewrite the above eq. by duality and minimax:

$$\sup_{\tilde{\pi} \in \Pi} \inf \left\{ \int_S \int_A A^\pi(s, a) d\tilde{\pi}(a|s) d\rho_\pi(s) + \lambda \left(\epsilon - \int_S C(\pi(.|s), \tilde{\pi}(.|s)) d\rho(s) \right) \right\}$$

s.t. $\lambda > 0$

\leq

$$\inf_{\lambda > 0} \left\{ \lambda \epsilon + \sup_{\tilde{\pi} \in \Pi} \left\{ \int_S \int_A A^\pi(s, a) d\tilde{\pi}(a|s) - \lambda C(\pi(.|s), \tilde{\pi}(.|s)) d\rho(s) \right\} \right\}$$

\leq

* $\inf_{\lambda > 0} \left\{ \lambda \epsilon + \int_S \sup_{\tilde{\pi} \in \Pi} \left\{ \int_A A^\pi(s, a) d\tilde{\pi}(a|s) - \lambda C(\pi(.|s), \tilde{\pi}(.|s)) \right\} d\rho(s) \right\}$

*

$$*: \lambda > 0 : (\text{if } \lambda = 0, * = \sup_{\tilde{\pi} \in \Pi} \int_A A^{\tilde{\pi}}(s, a) d\tilde{\pi}(a|s))$$

Kantorovich Duality: ($\psi(\cdot) = \frac{A^\pi(s, \cdot)}{\lambda}$ & $\phi(\cdot) = \inf_{a \in A} \{c(\cdot, a) - \psi(a)\})$

$$* = \sup_{\pi \in \Pi} \left\{ \int_A A^\pi(s, a) d\pi(a|s) - \lambda \sup_{\phi + \psi \in C} \left\{ \int_A \phi(a) d\pi + \int_A \psi(a) d\pi \right\} \right\}$$

$$\phi(a) + \psi(a') = \inf_{a \in A} \{c(a, a) - \psi(a)\} + \psi(a') \leq c(a, a')$$

$$\Rightarrow * \leq \sup_{\pi \in \Pi} \left\{ \int_A -\inf_{a' \in A} \{ \lambda c(a, a') - A^\pi(s, a') \} d\pi(a|s) \right\}$$

$$= \int_A \sup_{a' \in A} \{ A^\pi(s, a') - \lambda c(a, a') \} d\pi(a|s) \quad \text{***}$$

*** in
Dual Problem

\star

$$\inf \{ \lambda \in + \int_A \sup_{a' \in A} \{ A^\pi(s, a') - \lambda c(a, a') \} d\pi(a|s) d\rho_\pi(s)$$

s.t. $\lambda > 0$