

Name: Alfred Cueva  
Student No.: 2020-1783

## Exercise 7

### Data Science & Reinforcement Learning

Spring 2023

1. When is a policy gradient a better choice than a value-based method? When is a value-based method more preferred than a policy gradient? Please justify your answer.

- Policy gradient better when we have stochastic policies (optimal) <sup>And/or</sup> continuous action and smoother changes.  
- Value-based better for deterministic policies and/or action-space discrete.  
Computational efficient, stable

2. Explain what benefits Trust Region Policy Optimization (TRPO) has over simple linesearch-based policy gradient methods.

because <sup>in line searches</sup> learning rate should be estimated by simple linesearch-based tuning.  
because we don't overshoot we estimating the function, we use the surrogate function and the TR about the policy (always improves <sup>in TR</sup> and smooth changes in  $f$ ). Also more convenient when the variables are too large, we don't need to evaluate the constraint for all actions but rather focus on the trust region. Policy dist. restricted to change rapidly too. Always ~~is~~ stable too.  
<sub>and robust</sub>

3. Why is REINFORCE with Baseline not an actor-critic method?

Because ~~its~~ its state-value function is only used as baseline, not critic. State-value function isn't used for updating values of a state with estimates of subsequent states. Baseline doesn't play a role in policy update.

4. The policy gradient theorem for the episodic case establishes that

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_{\pi}(s, a) \nabla \pi(a|s, \theta).$$

Prove the above statement.

$$\nabla V_{\pi}(s) = \nabla \left[ \sum_a \pi(a|s) q_{\pi}(s, a) \right] \text{ for all } s \in \mathcal{S}.$$

$$\begin{aligned} &= \sum_a \left[ \nabla \pi(a|s) q_{\pi}(s, a) + \pi(a|s) \nabla q_{\pi}(s, a) \right] \\ &= \sum_a \left[ \nabla \pi(a|s) q_{\pi}(s, a) + \pi(a|s) \nabla \sum_{s', r} P(s', r|s, a) (r + V_{\pi}(s')) \right] \\ &= \sum_a \left[ \nabla \pi(a|s) q_{\pi}(s, a) + \pi(a|s) \sum_{s'} P(s'|s, a) \nabla V_{\pi}(s') \right] \end{aligned}$$

$$\begin{aligned} &= \sum_a \left[ \nabla \pi(a|s) q_{\pi}(s, a) + \pi(a|s) \sum_{s'} P(s'|s, a) \sum_a \left[ \nabla \pi(a|s') q_{\pi}(s', a) + \pi(a|s') \sum_{s'', r} P(s'', r|s', a) (r + V_{\pi}(s'')) \nabla V_{\pi}(s'') \right] \right] \\ &= \sum_{x \in \mathcal{S}} \sum_{k=0}^{\infty} P(s \rightarrow x, k, \pi) \sum_a \nabla \pi(a|x) q_{\pi}(x, a). \end{aligned}$$

$$\begin{aligned} \text{Then } \nabla J(\theta) &= \nabla V_{\pi}(s_0) = \sum_s \left( \sum_{k=0}^{\infty} P(s_0 \rightarrow s, k, \pi) \right) \sum_a \nabla \pi(a|s) q_{\pi}(s, a) = \sum_s \eta(s) \sum_a \nabla \pi(a|s) q_{\pi}(s, a) \\ &= \sum_{s'} \eta(s') \sum_s \frac{\eta(s)}{\sum_{s'} \eta(s')} \sum_a \nabla \pi(a|s) q_{\pi}(s, a) = \sum_{s'} \eta(s') \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_{\pi}(s, a) \\ &\propto \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_{\pi}(s, a). \end{aligned}$$