



Northeastern University  
CS 4180/5180 – Reinforcement Learning and Decision Making  
Fall 2019, Chris Amato

## Policy Gradients Assignment

Name: \_\_\_\_\_

Problem	Points
LINEAR FUNCTION APPROXIMATORS	/20
MISSING PSEUDOCODE	/20
SEMI-GRADIENT ONE-STEP EXPECTED SARSA	/20
GENERALIZE REINFORCE	/20
ELIGIBILITY VECTOR FOR SOFTMAX POLICY	/20
ELIGIBILITY VECTOR FOR GAUSSIAN POLICY	/20
<b>Total</b>	/120

### Instructions

- Don't cheat, write clearly, and staple your sheets!
- Questions about this assignment should be posted on Piazza for all to see.
- Contribute to the Piazza discussions, without providing direct answers.
- Graded out of 120 points for graduate students and 110 for undergrads.

**(20 pts.) LINEAR FUNCTION APPROXIMATORS**

Exercise 9.1 in the SB textbook. Show that tabular methods such as presented in Part I of this book are a special case of linear function approx. What would the feature vectors be?

Exercise 9.3 in the SB textbook. What  $n$  and  $c_{i,j}$  produce the feature vectors  $\mathbf{x}(s) = (1, s_1, s_2, s_1 s_2, s_1^2, s_2^2, s_1 s_2^2, s_1^2 s_2, s_1^2 s_2^2)^\top$ ?

**(20 pts.) MISSING PSEUDOCODE**

Exercise 10.1 in the SB textbook. We have not explicitly considered or given pseudocode for any Monte Carlo methods in this chapter. What would they be like? Why is it reasonable not to give pseudocode for them? How would they perform on the Mountain Car task?

**(20 pts.) SEMI-GRADIENT ONE-STEP EXPECTED SARSA**

Exercise 10.2 in the SB textbook. Give pseudocode for semi-gradient one-step *Expected* Sarsa for control.

**(20 pts.) GENERALIZE REINFORCE**

Exercise 13.2 in the SB textbook. Generalize the box on page 199, the policy gradient theorem (13.5), the proof of the policy gradient theorem (page 325), and the steps leading to the REINFORCE update equation (13.8), so that (13.8) ends up with a factor of  $\gamma^t$  and thus aligns with the general algorithm given in the pseudocode.

**(20 pts.) ELIGIBILITY VECTOR FOR SOFTMAX POLICY**

Exercise 13.3 in the SB textbook. In Section 13.1 we considered policy parameterizations using the softmax in action preferences (13.2) with linear action preferences (13.3). For this parameterization, prove that the eligibility vector is

$$\nabla \ln \pi(a \mid s, \boldsymbol{\theta}) = \mathbf{x}(s, a) - \sum_b \pi(b \mid s, \boldsymbol{\theta}) \mathbf{x}(s, b),$$

using the definitions and elementary calculus.

**(20 pts.)** ELIGIBILITY VECTOR FOR GAUSSIAN POLICY

Exercise 13.4 in the SB textbook. Show that for the gaussian policy parameterization (13.19) the eligibility vector has the following two parts:

$$\begin{aligned}\nabla \ln \pi(a \mid s, \boldsymbol{\theta}_\mu) &= \frac{\nabla \pi(a \mid s, \boldsymbol{\theta}_\mu)}{\pi(a \mid s, \boldsymbol{\theta})} = \frac{1}{\sigma(s, \boldsymbol{\theta})^2} (a - \mu(s, \boldsymbol{\theta})) \mathbf{x}_\mu(s) \text{ and} \\ \nabla \ln \pi(a \mid s, \boldsymbol{\theta}_\sigma) &= \frac{\nabla \pi(a \mid s, \boldsymbol{\theta}_\sigma)}{\pi(a \mid s, \boldsymbol{\theta})} = \left( \frac{(a - \mu(s, \boldsymbol{\theta}))^2}{\sigma(s, \boldsymbol{\theta})^2} - 1 \right) \mathbf{x}_\sigma(s).\end{aligned}$$