## Problem 3 - Assignment 16

Let us use the Softmax function to approximate the policy function

$$\pi(s,a;\theta) = \frac{e^{\phi(s,a)^T \theta}}{\sum_{b \in A} e^{\phi(s,b)^T \theta}}$$

Then $\log \pi(s,a;\theta) = \phi(s,a)^T \theta - \log \sum_{b \in A} e^{\phi(s,b)^T \theta}$

And the partial derivatives w.r.t. $\theta_i$ are given by:

$$\frac{\partial \log \pi(s,a;\theta)}{\partial \theta_i} = \phi_i(s,a) - \frac{\sum_{b \in A} \phi_i(s,b) e^{\phi(s,b)^T \theta}}{\sum_{b \in A} e^{\phi(s,b)^T \theta}}$$

$$= \phi_i(s,a) - \sum_{b \in A} \left[ \frac{e^{\phi(s,b)^T \theta}}{\sum_{b \in A} e^{\phi(s,b)^T \theta}} \right] \phi_i(s,b)$$

$$= \phi_i(s,a) - \sum_{b \in A} \pi(s,b;\theta) \phi_i(s,b)$$

$$= \phi_i(s,a) - E_\pi[\phi_i(s,\cdot)]$$

Therefore

$$\nabla_\theta \log \pi(s,a;\theta) = \phi(s,a) - E_\pi[\phi(s,\cdot)]$$

A simple way to enable the Compatible Function Approximation

$\nabla_w Q(s,a;w) = \nabla_\theta \log \pi(s,a;\theta)$ is to set $Q(s,a;w)$ to be linear in its features. We let the features of $Q(s,a;w)$ be $\nabla_\theta \log \pi(s,a;\theta)$, so we get:

$$Q(s,a;w) = w^T \nabla_\theta \log \pi(s,a;\theta)$$

Thus, it is easily observed that the required condition holds, since:

$$\nabla_w (w^T \nabla_\theta \log \pi(s,a;\theta)) = \nabla_\theta \log \pi(s,a;\theta)$$

$$E_\pi[Q(s,a;w)] = \sum_{a \in A} \pi(s,a;\theta) Q(s,a;w)$$

$$= \sum_{a \in A} \pi(s,a;\theta) w^T \nabla_\theta \log \pi(s,a;\theta)$$

$$= \sum_{a \in A} \pi(s,a;\theta) \sum_{i=1}^{m} w_i \frac{\partial \log \pi(s,a;\theta)}{\partial \theta_i}$$

$$= \sum_{a \in A} \pi(s,a;\theta) \sum_{i=1}^{m} w_i \frac{1}{\pi(s,a;\theta)} \frac{\partial \pi(s,a;\theta)}{\partial \theta_i}$$

$$= \sum_{a \in A} \sum_{i=1}^{m} w_i \frac{\partial \pi(s,a;\theta)}{\partial \theta_i}$$

$$= \sum_{i=1}^{m} w_i \frac{\partial}{\partial \theta_i} \left( \sum_{a \in A} \pi(s,a;\theta) \right) = \sum_{i=1}^{m} w_i \cdot \frac{\partial}{\partial \theta_i}(1) = 0$$