

Assignment 16

Miao-Chin Yen

March 8, 2022

Problem 3

Assume we have a finite action space \mathcal{A} . Let $\phi(s, a) = (\phi_1(s, a), \phi_2(s, a), \dots, \phi_m(s, a))$ be the features vector for any $s \in \mathcal{N}, a \in \mathcal{A}$. Let $\theta = (\theta_1, \theta_2, \dots, \theta_m)$ be an m -vector of parameters. Let the action probabilities conditional on a given state s and given parameter vector θ be defined by the softmax function on the linear combination of features: $\phi(s, a)^T \cdot \theta$, i.e.,

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

- Evaluate the score function $\nabla_{\theta} \log \pi(s, a; \theta)$
- Construct the Action-Value function approximation $Q(s, a; \mathbf{w})$ so that the following key constraint of the Compatible Function Approximation Theorem (for Policy Gradient) is satisfied:

$$\nabla_{\mathbf{w}} Q(s, a; \mathbf{w}) = \nabla_{\theta} \log \pi(s, a; \theta)$$

where \mathbf{w} defines the parameters of the function approximation of the Action-Value function.

- Show that $Q(s, a; \mathbf{w})$ has zero mean for any state s , i.e. show that

$$\mathbb{E}_{\pi}[Q(s, a; \mathbf{w})] \text{ defined as } \sum_{a \in \mathcal{A}} \pi(s, a; \theta) \cdot Q(s, a; \mathbf{w}) = 0 \text{ for all } s \in \mathcal{N}$$

Answer:

$$\begin{aligned} \log \pi(s, a; \theta) &= \theta \cdot \phi(s, a)^T - \log\left(\sum_{b \in \mathcal{A}} e^{\phi(s, b)^T \cdot \theta}\right) \\ \frac{\partial \log \pi(s, a; \theta)}{\partial \theta_i} &= \phi_i(s, a) - \frac{\sum_{b \in \mathcal{A}} \phi_i(s, b) \cdot e^{\phi(s, b)^T \cdot \theta}}{\sum_{b \in \mathcal{A}} e^{\phi(s, b)^T \cdot \theta}} \\ &= \phi_i(s, a) - \sum_{b \in \mathcal{A}} \frac{e^{\phi(s, b)^T \cdot \theta}}{\sum_{b \in \mathcal{A}} e^{\phi(s, b)^T \cdot \theta}} \cdot \phi_i(s, b) \\ &= \phi_i(s, a) - \sum_{b \in \mathcal{A}} \pi(s, b; \theta) \cdot \phi_i(s, b) \\ &= \phi_i(s, a) - \mathbb{E}_{\pi}[\phi_i(s, \cdot)] \\ &\implies \nabla_{\theta} \log \pi(s, a; \theta) = \phi(s, a) - \mathbb{E}_{\pi}[\phi(s, \cdot)] \end{aligned}$$

Construct the Action-Value function approximation as follows:

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

Then we can satisfy the key constraint of the Compatible Function Approximation Theorem

$$\nabla_{\mathbf{w}} Q(s, a; \mathbf{w}) = \nabla_{\theta} \log \pi(s, a; \theta)$$

And,

$$\begin{aligned}\sum_{a \in \mathcal{A}} \pi(s, a; \boldsymbol{\theta}) \cdot Q(s, a; \boldsymbol{w}) &= \sum_{a \in \mathcal{A}} \pi(s, a; \boldsymbol{\theta}) \cdot \boldsymbol{w}^T \cdot \nabla_{\theta} \log \pi(s, a, \boldsymbol{\theta}) \\ &= \sum_{a \in \mathcal{A}} \boldsymbol{w}^T \cdot \nabla_{\theta} \pi(s, a, \boldsymbol{\theta}) \\ &= \boldsymbol{w}^T \cdot \nabla_{\theta} \left(\sum_{a \in \mathcal{A}} \pi(s, a, \boldsymbol{\theta}) \right) \\ &= \boldsymbol{w}^T \cdot \nabla_{\theta} 1 \\ &= \boldsymbol{w}^T \cdot \mathbf{0} \\ &= 0\end{aligned}$$