

Policy Gradient Methods

Tabular methods and Function Approximators are value-based which means the policy looks at the value of the state-actions Q and choose the action with optimal estimated return. In policy gradient methods, we use the function approximators to not predict the Q values but to estimate the probabilities of taking each action:

$$\pi(a|s, \theta) \in [0, 1] \quad (1)$$

and thus the function approximator **is** the policy that gives out the distribution of probabilities over the actions rather than the values for you to read:

$$\pi(s, \theta) = [p(a_1), p(a_2), \dots, p(a_n)] \quad (2)$$

and your action to take, a , is sampled from this distribution:

$$a \sim \pi(a|s, \theta) \quad (3)$$

and thus your policy and actions change more smoothly during learning, whereas in value-based methods, you manually choose the action using greedy or ϵ -greedy methods:

$$a = \arg \max_a \hat{Q}(s, a). \quad (4)$$

Representing Policies using Neural Networks

Neural networks can learn the weights and propagate the extracted features of the state vector to approximate the values of the state:

$$\hat{y} = \phi_2(\phi_1(x \cdot w_1) \cdot w_2). \quad (5)$$

To train the neural network to output the vector of probability, we compress the output values to values between 0 and 1 such that:

$$\hat{y} = [0, 1]^n \quad (6)$$

where n is the number of possible actions to take.

Softmax Function

In order to achieve the compression, we apply softmax to normalize the values so that not only they compress to $[0, 1]$, they also all sum up to 1:

$$\sigma(x_i) = \frac{e^{x_i}}{\sum e^{x_i}} \quad (7)$$

where the values of the individual outputs are raised to powers of e to get a vector of probability distribution over the possible actions:

$$\pi(a|s, \theta) = [\sigma(a_1), \sigma(a_2), \dots, \sigma(a_n)]. \quad (8)$$

Policy Evaluation

We define performance measure of the policy/neural network as $J(\theta)$, so that we can compare them such that:

$$J_{\pi_1}(\theta) > J_{\pi_2}(\theta) \Rightarrow \text{We consider } J_{\pi_1}(\theta) \text{ is better than } J_{\pi_2}(\theta). \quad (9)$$

As such the optimal policy has the optimal parameters θ that give us the maximal performance:

$$\pi_*(a|s, \theta) = \arg \max_{\theta} J(\theta). \quad (10)$$

We obtain the performance estimate $\hat{J}(\theta)$ from experience samples during the simulations, and we perform the Stochastic Gradient Ascent to improve the performance with regards to its parameters:

$$\theta_{t+1} = \theta_t + \alpha \nabla J(\theta) \quad (11)$$

where

$$\nabla J(\theta) = \left[\frac{\partial \hat{J}(\theta)}{\partial \theta_1}, \frac{\partial \hat{J}(\theta)}{\partial \theta_2}, \dots, \frac{\partial \hat{J}(\theta)}{\partial \theta_n} \right]. \quad (12)$$

Policy Gradient Theorem