

Plicy Gradient methods in the Evolutionary Pricing Game

June 16, 2023

Abstract

The deep reinforcement learning model of an agent in duopoly multi-round pricing game and applying policy gradient methods to find the optimal strategy.

1 Policy Gradient methods

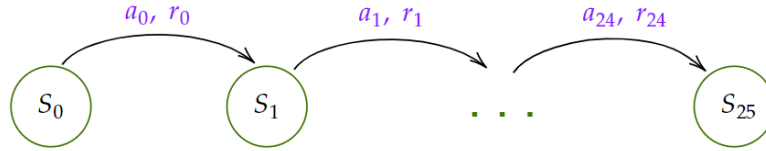


Figure 1 environment model

f1equi

In *action-value* methods such as Q-learning, the value of each state or state-action pair is learned and based on these values, the policy would be determined. However, in *Policy Gradient* methods, the policy is learnt directly without using value functions for action selection.

Since the number of states in our model is large, we need a function approximator to parametrize the action preferences. We use an artificial neural network for this purpose and we show the vector of connection weights in our network as θ .

The actions at each state are chosen in a way that action with higher valuation is more likely to be chosen.

1.1 Action space

The reward in each stage of the multi-round pricing game is determined by $(D_i - P) \times (P - C)$ in which D_i is the agent's demand potential at the start of stage i and C is the agent's product cost. $P^* = \frac{D_i + C}{2}$ maximises this quadratic concave function; we refer to P^* as *myopic price* or *monopoly price* at stage i . Playing the myopic price results in maximum payoff at the current stage but the demand potential for the next round would be affected, meaning it does **not** necessarily result in the maximum overall return.

We considered actions as the value below the monopoly price that should be played, in order to attract more demand and consequently more reward in later stages. In our model the action-space is discrete. In each stage, action $a \in \{0, 1, \dots, 19\}$ with $\text{step} = 3$ that means the price can be 0, 3, 6, ..., or 57 units less than the stage's myopic price.

For example, for the low-cost player ($C = 57$), if $D_i = 180$ at the start of stage i and action $a = 5$:

$$P^* = \frac{180 + 57}{2} = 118.5 \rightarrow P = P^* - (a \times \text{step}) = 118.5 - 5 \times 3 = 103.5$$

2 Monte Carlo Policy Gradient

3 Reinforcement with Baseline

4 Actor-Critic method