
Algorithm 1 Learning Algorithm

```
Randomly initialize critic network  $V_\pi^\phi(s_{it}, S_t)$  and actor network  $a^\theta(s_{it}, S_t, q_t)$  with
parameters  $\phi$  and  $\theta$ 
Initialize state  $S_t, \{s_{it}\}_{i \in I}$  according to user-supplied values/distributions
if PRETRAIN_CHOICE

    Using user-provided initial values for choices as labels, train the policy
    network to select those actions given initial state
    # This is just an attempt to initialize the policy network better than randomly

end
# Burn-in: Approach the stochastic steady state given the initial policy.
for t = 1:T_burn

    Guess values for equilibrium variables  $q_t$  (e.g. prices)
    while error <  $\delta$ 
        Sample actions  $a_{it} =$ 
         $a^\theta(s_{it}, S_t, q_t)$  for all agents  $i$ 
        Compute error  $\leftarrow$  equilibrium_conditions( $\{a_{it}, s_{it}\}_{i \in I}, S_t, q_t$ )
        (e.g. market clearing conditions)
    end
    Draw shocks  $\{\varepsilon_{it}\}_{i \in I}, \varepsilon_t$ 
    Compute next state  $(\{s_{it+1}\}_{i \in I}, S_{t+1}) =$ 
     $T(\{a_{it}, \varepsilon_{it}, s_{it}\}_{i \in I}, S_t, q_t, \varepsilon_t)$ 
end
# Estimate the value function given the initial policy
for t = 1:T_trainvalue

    Solve for equilibrium as above
    for j = 1:M
        Draw shocks  $\{\varepsilon_{it}\}_{i \in I}, \varepsilon_t$ 
        Compute next state  $\left(\{s_{it+1}^j\}_{i \in I}, S_{t+1}^j\right) =$ 
         $T(\{a_{it}, \varepsilon_{it}, s_{it}\}_{i \in I}, S_t, q_t, \varepsilon_t)$ 
        Compute for each  $i \in I$  the TD target  $y_{it}^j =$ 
         $r_i + \gamma V_\pi^\phi(s_{it+1}^j, S_{t+1}^j)$ 
    end
    Estimate the expected average TD error gradient  $\nabla_\phi \delta_t^\phi =$ 
     $\nabla_\phi \frac{1}{N} \sum_{i=1}^N \left( \frac{1}{M} \sum_{j=1}^M y_{it}^j - V_\pi^\phi(s_{it}, S_t) \right)^2$ 
    Update  $\phi$  according to the estimated gradient  $\nabla_\phi \delta_t^\phi$ 
end
```

Algorithm 2

```
# Train the policy while moving toward the equilibrium steady state
and updating the value function
for episode = 1:M
    for t = 1:T_trainpolicy
        solve for equilibrium as above
        for j = 1:M
            Draw
            Estimate  $\nabla_{\theta} \mathbb{E}[r_t + V(s_{t+1}, S_{t+1}) \mid S_t] \approx$ 
 $\frac{1}{N} \sum_{i \in I} [\nabla_{\theta} (r_{it} + V_{\pi}^{\phi}(s_{it+1}, S_t))]$ 

            where  $\nabla_{\theta} (r_{it} + V_{\pi}^{\phi}(s_{it+1}, S_t))$  is computed directly using autodiff.
            Add this gradient to the policy gradient buffer.
            For discrete actions, update  $\theta$ 
            Compute for each  $i \in I$  the TD target  $y_{it} =$ 
 $r_i + \gamma V_{\pi}^{\phi}(s_{it+1}, S_{t+1})$ 
            Estimate the expected TD error  $\delta_t =$ 
 $\mathbb{E}[\delta_{it} \mid S_t] \approx \frac{1}{N} \sum_{i \in I} (y_{it} - V_{\pi}^{\phi}(s_{it}, S_t))^2$ 
            Compute the gradient  $\nabla_{\phi} \frac{1}{N} \sum_{i \in I} (y_{it} - V_{\pi}^{\phi}(s_{it}, S_t))^2$  and add it to
            the TD error gradient buffer
        end
        Update  $\theta$  according to the policy gradient buffer
        Update  $\phi$  according to the TD error gradient buffer
        Update state so that the first state of the next episode is
        the last state from the previous episode
    end
end
```
