

Advanced Machine Learning: Assignment 2

Seung-Hoon Na

December 03, 2017

1 Policy gradient theorem

In the episodic case, suppose that every episode starts in some particular state s_0 and the objective (performance) $J(\theta)$ is defined as:

$$J(\theta) = v_{\pi_\theta}(s_0)$$

where $v_{\pi_\theta}(s_0)$ is the true value function for π_θ , the policy determined by θ .

The *policy gradient theorem* is stated as:

Theorem 1 (Policy gradient theorem)

$$\nabla J(\theta) = \mathbb{E}_\pi \left[\sum_a q_\pi(S_t, a) \nabla_\theta \log \pi(a|S_t, \theta) \right]$$

Prove the policy gradient theorem. (refer to section 13.2 in Sutton's textbook).

Hint: Take the following steps.

1. Derive the following equation.

$$\nabla J(\theta) = \nabla v_\pi(s_0) \propto \sum_s d^\pi(s) \sum_a \nabla \pi(a|s) q_\pi(s, a) \quad (1)$$

where $\eta(a)$ denotes the number of time steps spent, on average, in state s in a single episode, which is formulated as (refer to Section 9.2 in Sutton's book):

$$\eta(s) = h(s) + \sum_{s'} \eta(s') \sum_a \pi(a|s') p(s|s', a)$$

and $d^\pi(s)$ is the stationary distribution of Markov chain under π_θ , being formulated as:

$$d^\pi(s) = \frac{\eta(s)}{\sum_{s'} \eta(s')}$$

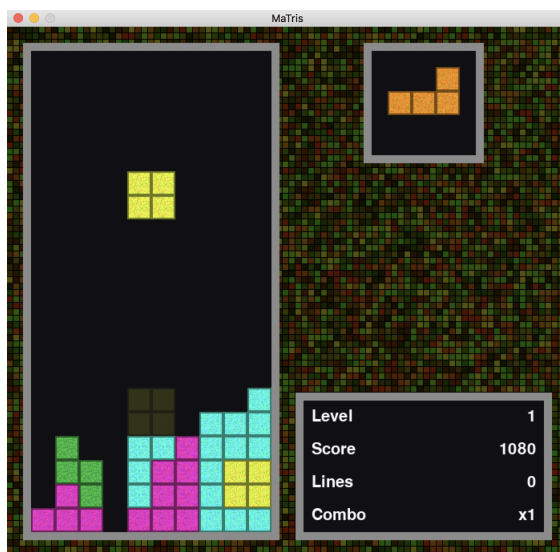
2. Obtain the final theorem by rewriting Eq. 1 in terms of expectation. (Provide the detailed derivation how the summation part is changed using the expectation).

2 Deep reinforcement learning: DQN and Actor-Critic

In this problem, you should implement DQN and actor-critic model based on convolutional networks and apply them to learn actions for Tetris game.

Refer to the following pygame MaTris codes:

<https://github.com/SmartViking/MaTris> which need to be extended for Here is the screen shot when playing the MaTris game.



2.1 DQN (Deep Q-network: implementation)

DQN takes the replay memory \mathcal{D} consisting of experience replays (s, a, r, s') , randomly samples a mini-batch \mathcal{D}_i from \mathcal{D} , and updates the parameters of action-value networks to reduce the following loss function (squared error).

$$L_i = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} [r + \gamma \max_{a'} Q(s', a'; \mathbf{w}_{old}) - Q(s, a, \mathbf{w})]^2$$

Here is the detailed algorithm presented in DQN paper by DeepMind. (<https://www.cs.toronto.edu/~vmnih/docs/dqn.pdf>)

Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N
Initialize action-value function Q with random weights
for episode = 1, M **do**
 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$
 for $t = 1, T$ **do**
 With probability ϵ select a random action a_t
 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$
 Execute action a_t in emulator and observe reward r_t and image x_{t+1}
 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D}
 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D}
 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$
 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3
 end for
end for

Implement a generic DQN that learns $Q(s, a, w)$, which is designed based on deep convolutional networks (using tensorflow), and apply the implemented DQN to learn Tetris world.

Refer to the DeepMind's DQN paper to see how $Q(s, a, w)$ is designed using convolutional networks.

Note: The codes for DQN need to have the high-level of modularity such that they can be easily extended for other games and RL worlds.

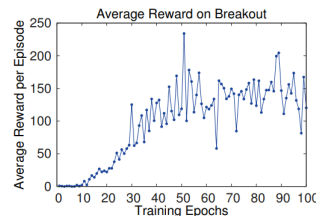
You may need to revise the MaTris codes to extract the rewards and the state images that correspond to he current screen.

Here are additional comments:

1. Your implementation need to be compatible with tensorflow 1.0.
2. You can fix the discount factor γ to a value in the range of $[0.7, 0.95]$.
3. Provide a short README file for executing the implemented codes.

2.2 DQN (Deep Q-network): Performance curves of DQN in Tetris world

Draw a performance curve which shows how average rewards are changed according to the number of the training epochs, similar to the following curves presented in the DeepMind's DQN paper:



2.3 DQN (Deep Q-network): Simulation of playing Tetris performed by the learned agent

Write the simulation code that automatically plays MaTris by the learned policies. Once the simulation starts, an agent automatically play Tetris using his learned policy on the screen.

Of course, at the training codes, you first need to store the learned model from DQN in a separate file, named a *DQN model file*.

The simulation code loads a DQN model file and uses it as policy to automatically play Tetris.

2.4 Actor-critic model: Implementation

Actor-critic model consists of two networks – a value network $v(S, \mathbf{w})$ and a policy network $\pi(S, a, \theta)$. To learn parameters of a value network, TD-learning with the function approximation is applied. To learn parameters of policy network, use the following modified policy gradient that uses advantage function $A^{\pi\theta}(s, a)$ instead of original q-value:

$$\nabla J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla \log \pi(a|S_t, \theta) A^{\pi\theta}(S, a)]$$

The algorithm for actor-critic model is given below (13.5 in Sutton’s textbook).

One-step Actor–Critic (episodic)

Input: a differentiable policy parameterization $\pi(a|s, \theta)$

Input: a differentiable state-value parameterization $\hat{v}(s, \mathbf{w})$

Parameters: step sizes $\alpha^{\theta} > 0$, $\alpha^{\mathbf{w}} > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $\mathbf{w} \in \mathbb{R}^d$

Repeat forever:

 Initialize S (first state of episode)

$I \leftarrow 1$

 While S is not terminal:

$A \sim \pi(\cdot|S, \theta)$

 Take action A , observe S', R

$\delta \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$ (if S' is terminal, then $\hat{v}(S', \mathbf{w}) \doteq 0$)

$\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} I \delta \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w})$

$\theta \leftarrow \theta + \alpha^{\theta} I \delta \nabla_{\theta} \ln \pi(A|S, \theta)$

$I \leftarrow \gamma I$

$S \leftarrow S'$

Implement a generic actor-critic model that learns $v(S, \mathbf{w})$ and $\pi(s, a, \mathbf{w})$, which are designed based on deep convolutional networks

(using tensorflow), and apply the implemented actor-critic model to learn our Tetris world.

Likewise, the codes for actor-critic model need to have the high-level of modularity such that they can be easily extended for other games and RL worlds.

The MaTris code needs to be modified and extended to extract state images and rewards.

All other settings (e.g. discount factor, tensorflow version, minibatch size) should be similar to those of DQN

2.5 Actor-critic model: Application and simulation

1. **Draw a performance curve of actor-critic model which shows how average rewards are changed according to the number of the training epochs, similar to the performance curves of the DeepMind's DQN paper:**
2. **Write the simulation code of actor-critic model that automatically plays MaTris by the learned policies, and do the simulation on the game screen.**

Notes: You need to prepare README files that describe how to run your codes – both for printing performance curve and simulating the learned model.

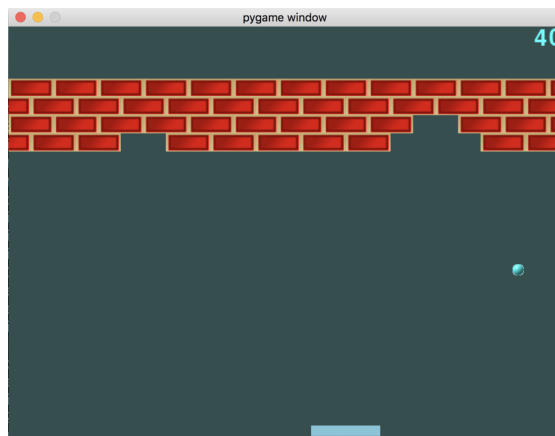
3 DQN and Actor-Critic for Breakout game

Extend your implementation of DQN and AC (actor-critic) codes and apply to Breakout game.

The following is the pygame code that you need to extend or revise:

<https://github.com/johncheetham/breakout>

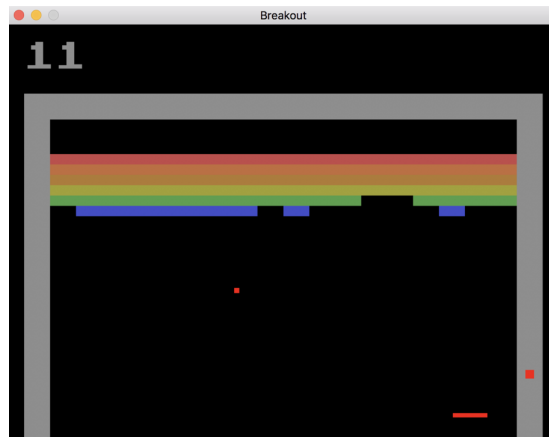
Here is the screen shot when playing the breakout game.



Or, you can revise and extend the following pygame code that implements a Atari-style breakout:

<https://github.com/aknuck/Atari-Breakout>

Here is the screen shot when playing the breakout game.



Likewise, you need to submit the following codes and corresponding README, with a short report.

1. DQN for controlling Breakout: Training (and store the learned model in a separate file)
2. DQN for controlling Breakout: Testing – performance curve
3. DQN for controlling Breakout: Testing – simulation (auto-play by a game agent)
4. AC for controlling Breakout: Training (and store the learned model in a separate file)
5. AC for controlling Breakout: Testing – performance curve
6. AC for controlling Breakout: Testing – simulation (auto-play by a game agent)