

CSC411 Machine Learning

Project 4: Reinforcement Learning

Ariel Kelman

Student No: 1000561368

Gideon Blinick

Student No: 999763000

3 April 2018

1 Introduction

1.1 Results & Report Reproducibility

All results and plots can be generated using the code in `tictactoe.py`. All code is python 3, run with Anaconda 3.6.3. Running the code will save all generated images in the `resources` folder, where they are used by \LaTeX . Note that some of the sections require the code for other sections (in the same file) to be run first. To reproduce the report, simply run it through \LaTeX . This will pull the most recently generated figures from the `resources` folder.

1.2 Tic Tac Toe Environment

The `Environment` class provides the functionality required to play a game of tic tac toe, including against an opponent who plays a random legal move. The game is represented by the `grid` attribute, which is a `numpy` array representing the 9 positions. Each position can have a 0, 1 or 2, representing an empty position, or one filled by player 1 (*X*) or 2 (*O*) respectively. The `turn` attribute can have a value of 1 or 2, and represents which player has the next move. The `done` attribute is a boolean that indicates whether a game has been completed (either with a win, or when the board is full).

The `step()` and `render()` methods allow a tic tac toe game to be played and displayed with text output. Using these methods, a game was played, resulting in the following output:

```
...  
.x.  
...  
=====  
o..  
.x.  
...  
=====  
o.x  
.x.  
...  
=====  
o.x  
ox.
```

```

...
=====
O.X
OX.
X..
=====

```

2 Policy

The `Policy` class implements a neural network that learns to play tic tac toe. The provided starter code was modified to be a one-hidden layer neural network; the final code is shown in the code block below.

```

class Policy(nn.Module):
    """
    The Tic-Tac-Toe Policy
    """
    def __init__(self, input_size=27, hidden_size=64, output_size=9):
        super().__init__()
        self.Linear1 = nn.Linear(input_size, hidden_size)
        self.Linear2 = nn.Linear(hidden_size, output_size)

    def forward(self, x):
        h = F.relu(self.Linear1(x))
        out = F.softmax(self.Linear2(h))
        return out

```

In choosing an action, the state is represented as a 27-dimensional vector, using a one-hot encoding type scheme. The first nine elements are 1 if the corresponding location (moving horizontally, and then from top to bottom), and 0 otherwise. Similarly, the next 9 elements are 1 if an *X* (representing player 1) is in the corresponding location; while the last nine elements provide the same functionality for *O* (player 2).

The policy outputs a nine-dimensional vector, which is a probability distribution that is sampled to choose the move for the policy. Thus the policy is stochastic - the `select_action()` function samples this distribution to choose a move.

3 Policy Gradient

The `compute_returns()` function computes the returns based on the reward at the end of a game. The following code block shows how the returns are calculated.

```

def compute_returns(rewards, gamma=1.0):
    """
    Compute returns for each time step, given the rewards
    """
    k = len(rewards)
    rewards = np.array(rewards)
    gammas = np.array([gamma**(i) for i in range(k)])
    G = [sum(rewards[i:]*gammas[:k-i]) for i in range(k)]
    return G

```

The weights are updated on the conclusion of a game, though this means that the policy does not improve during a game (this is a minor cost considering how short the games are). For updates to the policy to have any meaning in the middle of the game, there would need to be rewards for intermediate

states. This would require having a metric that evaluates a tic tac toe position. While this could be done in tic tac toe, it would be quite difficult in more complicated games, where it is not obvious how a position should be rewarded. Having rewards only at the end still improves the policy choices earlier in the game (if an early position leads to a win, that position will be weighted more positively when propagating backwards). Thus updating the weights at the end of an episode provides a simple way to provide feedback and improve the policy.

4 Rewards

5 Training

6 Win Rate

7 First Moves

8 Limitations