

# Homework 6 (due January 23, 2018) – Reinforcement learning

## I) warmup: Tic-Tac-Toe in Python

To get used to the main concept, we'll start with a simple example: Q-learning for Tic-Tac-Toe. This can be done without neural network, so we won't use tensorflow in this part.

Tip: There are many possible ways to implement Q-learning. Part of the homework is to find suitable ones yourself. Probably, the first you come up will not be the most suitable one, so be prepared to throw away (part of) your code at some time and try something else.

### I.1) Tic-Tac-Toe

If you don't remember it, familiarize yourself with the rules of Tic-Tac-Toe (<https://en.wikipedia.org/wiki/Tic-tac-toe>)

### I.2) Representations

Choose a representation of the playing board (9 locations, each either empty or filled with a player's mark). You should be able to

- access individual locations of the board by an index and read out or overwrite them,
- convert the board into a *hashable* type (typically *integer* or *string*)

Write a subroutine that checks if for any board configuration, if Player 1 has won (output +1), Player 2 has won (output -1). Otherwise, output 0.

Write a subroutine that checks for any board configuration if there's any legal moves (=actions) left to make.

Write a subroutine that for a given board configuration, move location, and player ID updates the board accordingly, and that computes the reward of the action (-1/0/+1 depending on the value of the new board configuration).

### I.3) Random play

Write a program that plays randomly. For a given number of rounds:

- create an empty board
- repeat
  - let Players 1 and 2 alternately make random (legal) moves
- until there are no moves left or one of the players has won the game

a) run for 10000 games (if this takes longer than a minute or two, better rethink your implementation).

b) evaluate in three ways: i) record sequence of all rewards and plot its cumulative sum as a curve,

ii) compute the *total reward* (i.e. the sum of the sequence), iii) build a histogram of how many wins/draws/losses there were for Player 1 won

iii) modify the code such that in each round it is randomly chosen if Player 1 or Player 2 move first.

iv) reevaluate as in b)

### I.4) Stronger opponents

a) Replace the random Player 2 by a *smart* one: in any situation, it checks if there's an immediate winning possibility available. If yes, it takes it. If not, it plays randomly. Let it run again a random player 1, what's the outcome?

b) Make the opponent even *smarter*: it should still make winning moves if it can, but if can't it should check if the opponent has a winning move and if yes, it blocks that. If not, it plays randomly.

c) Perform a tournament of all all options (random, smart, smarter) for Players 1 and Player 2 and record the *total reward* in each case.

To facilitate this, you might want to refactor your code such that the players' strategies can be switched without having to every time change a lot of code.

### I.5) Q-learning

We now integrate *Q-learning* for Player 1, to enable it to learn from its experience.

- The *state space* will be the set of all possible board configurations.
  - The *action space* is the 9 locations of the board where the player can make moves (i.e. make a symbol). The action space is the same for all states, even if the rules forbid some moves in some situations.
  - The *transitions* are: when Player 1 has made a move, that appears on the board. Then (unless the game is over) Player 2 makes a move according to some strategy. The new state is the board position after this.
  - The *reward* of an action is the quality of the position resulting from it (+1, if Player 1 creates three-in-a-row, -1 if Player 2 creates three-in-a-row afterwards, 0 in all other cases).
- a) Create a data structure that allows table lookups for the *Q* function. For any board *s* state it should output 9 values, one per action *a*.

Tip: a convenient option is a *defaultdict* from Python's *collections* package. It can be indexed by any hashable type and the *default* (e.g. in combination with a *lambda* function) allows receiving a default value (e.g. all zeros) even for entries that haven't been written yet.

- b) Write a routine `updateQ(s,a,s',r)` that performs a *Q-learning update*: after Player 1 in state *s* made action *a*, which resulted in a new state *s'* and a reward value *r*:

- $Q(s)[a] := r + \max_{a'} Q(s')[a']$  # the value for action *a* in state *s* is *reward* plus value of best continuation

Note: in comparison to the more general form  $Q(s)[a] := (1 - \alpha)Q(s)[a] + \alpha(r + \gamma \max_{a'} Q(s')[a'])$ , we don't use a discounting factor (i.e.  $\gamma = 1$ ) and the learning rate is  $\alpha = 1$ .

- c) Integrate the Q-learning update in the above random play routine by calling it every time Player 1 has made a move, and the new state and reward has been determined.

### I.6) Putting things together

For Player 1 to learn to play better over time, it must make use of the learned *Q*-function.

- a) Implement a routine that makes *greedy* moves for Player 1:
- in any state *s*, pick the action *a* with highest values of  $Q(s)[a]$  that is legally possible (i.e. the place on the board must be empty)

Integrate the greedy Player 1 into the code from I.5) and let it play against random, smart, and smarter Players 2.

- a) Compare the total rewards of a) to the results from I.4.c) for the 'smarter' Player 1. Is Q-learning doing better or worse?
- b) Now plot the curves of cumulative rewards. What do you observe? What would happen if you had played 20.000 instead of 10.000 games?
- c) (optional) Can you imagine a situation in which Q-learning with the *greedy* player does less well than it could? How could you prevent that?

## II) Tic-Tac-Toe in Tensorflow

We now switch from table-based Q-learning to using a neural network. Note that technically this isn't necessary for Tic-Tac-Toe, as the state space is small and the simpler table-based system works well enough. We just use it here to practice the tensorflow integration.

## II.1) Q-function in tensorflow

Instead of the above look-up table, implement a  $Q$  function (approximator) using a neural network.

- as input, feed the board configuration (in the original form) via a placeholder
- convert the board to a one-hot encoding (if it wasn't already) and flatten it
- create one densely connected layer with 10 hidden units and a nonlinearity of your choice
- create one densely connected layer with 9 output units (each output is the  $Q$ -values of one of the action)

For the fully connected layers, use an initialization with very small values, such that the  $Q$ -values don't differ too much initially.

## II.2) updating the $Q$ function

Updating the  $Q$  function is the main problem in deep  $Q$ -learning, as one cannot simply overwrite values by others as for the function table. Instead, we have to run a few steps of network training to make the specific output we want to influence closer to the value we want it to have. There's multiple ways to do this, here is a simple (though not very elegant) one:

- create a `tf.float32` placeholder  $t$  for the target value (the right hand side of the  $Q$ -update)
- create a `tf.int32` placeholder  $a$  for the index of the action that should be affected
- define a tensor that is the squared difference between the  $a$ -th output entry of the  $Q$ -function network and  $t$
- create an operation for minimizing this squared difference (e.g. using a `tf.train.GradientDescentOptimizer`. What's a good learning rate? Your guess is as good as mine...)

Write a new `updateQ(s,a,s',r)` routine that

- evaluates  $Q(s')$  (using `sess.run`) and writes it to a python vector  $q$
- computes the desired right hand side (in python),  $r + \max_{a'} q[a']$
- call the above minimization operation a certain number of steps (e.g. 10) to achieve the  $Q$ -update

(optional challenge): replace as many python/numpy operations by tensorflow operations as possible (e.g. computing  $\max_{a'} Q(s')[a']$ ). Try to find a way to replace the two `sess.run` calls by a single one. You might have to look up `tf.stop_gradient` for this.

## II.3) experiments

Let network-based  $Q$ -learner play against 'random', 'smart' and 'smarter' Player 2. Does it work better or worse than the table-based learning in part I)? Why do you think that is? Can you think of anything to improve the system (further)?

## III) Connect-4 (due January 30th)

Write a network-based system that learn to play Connect-4 ([https://en.wikipedia.org/wiki/Connect\\_Four](https://en.wikipedia.org/wiki/Connect_Four)) against a random opponent, or the analog of the 'smart' and 'smarter' opponent. Try to find a better architecture than the one in II). Maybe convolutional? What filter size makes sense?

## IV) (optional) Tic Tac Toe - 5 in Row

Write a network-based system that learns to play an extended version of Tic-Tac-Toe, in which the player have to achieve 5 in a row/column/diagonal on a 15x15 board. Note that you will have to come up with a reasonable opponent as well.

## Hand-in requirements

Upload your code to I), II) and eventually III), as well as the results table from I.4.c) and the curves of cumulative reward from I.6.b) to the IST *git* server.