Problem Set 3: Reinforcement Learning

The goal of this problem set is to implement and get familiar with Reinforcement Learning.

To run the autograder, type the following command in the terminal:

```
python autograder.py
```

If you wish to run a certain problem only (e.g. problem 1), type:

```
python autograder.py -q 1
```

where 1 is the number of the problem you wish to run.

Instructions

In the attached python files, you will find locations marked with:

```
#TODO: ADD YOUR CODE HERE
NotImplemented()
```

Remove the NotImplemented() call and write your solution to the problem. **DO NOT MODIFY ANY OTHER CODE**; The grading of the assignment will be automated and any code written outside the assigned locations will not be included during the grading process.

IMPORTANT: You must document your code (explain the algorithm you are implementing in your own words within the code) to get the full grade. Undocumented code will be penalized.

For this assignment, you should submit the following files only:

- value_iteration.py
- options.py
- reinforcement learning.py

Put the files in a compressed zip file named solution.zip which you should submit to Blackboard.

Problem Definitions

There is one environment defined in this problem set:

1. **Grid World**: where the environment is a 2D grid where the player can move in one of 4 directions (U, D, L, R). The actions are noisy so the end result may be moving along one of the 2 directions orthogonal to the desired direction. Some locations "#" are occupied with walls so the player cannot stand on them. Some other locations "T" are terminal states which ends the episode as soon as the player stands on them. Each location has an associated reward which is given to player if they do an action that gets them to be in that state in the next time step. The markov decision process and environment of the grid world is implemented in grid.py and the environment instances are included in the grids folder.

You can play a grid world game by running:

```
# For playing a grid (e.g. grid1.json)
python play.py grids\grid1.json
```

You can also let an learning agent play the game in your place (e.g. a Q-Learning Agent) as follow:

```
python play.py grids\grid1.json -a q learning -m models/model.json
```

NOTE: In addition to the agent, we supply a model file from which the agent will read its data (e.g. Q-values for Q-learing & SARSA agents, weights for approximate Q-learning and Utilities for value iteration agents). If we don't supply a model file, the agent will play using the initial values of their learnable parameters.

To train an agent and save its data in model file, use train.py as follows:

```
# For training a q_learning agent on grid1.json for 1000 iterations where each episode python play.py q_learning grids\grid1.json models/model.json -i 1000 -sl 100
```

The agent options are:

- human where the human play via the console
- random where the computer plays randomly
- value_iteration where the agent uses the learned utilities (via value iteration) and the MDP to decide on the action to take.
- sarsa where the agent uses the learned Q-value (via SARSA) to decide on the action to take.
- q_learning where the agent uses the learned Q-value (via Q-Learning) to decide on the action to take.
- q_learning_approx where the agent uses the learned weights (via Linear Approximate Q-Learning) to decide on the action to take.

To get detailed help messages, run play.py and train.py with the -h flag.

Important Notes

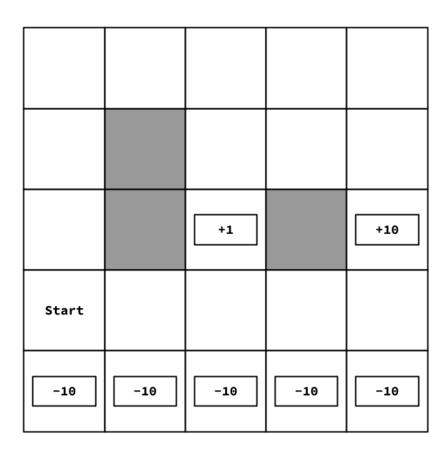
This problem set relies a lot on randomness and Reinforcement Learning usually does not converge to the same results when different random seeds are used. So it is essential to follow the instructions written in the comments around the TODOs to acheive the same results. In addition, while acting, if two actions have the same value (Q-value, expected utilities, etc.), pick the action that appears first in the list returned by mdp.get_actions or env.actions().

Problem 1: Value Iteration

Inside value_iteration.py, modify the functions marked by a **TODO** to complete the ValueIterationAgent class. Note that the Reward R(s, a, s') is a function of the current state, the action and the next state, so use the appropriate version of the bellman equation:

$$U(s) = \max_a \sum_{s'} P(s'|s,a)[R(s,a,s') + \gamma U(s')]$$

Problem 2: Parameter Selection



As shown in figure above, this MDP has 2 terminal states with positive rewards (one is close to the start with a reward of +1 and far from the start with a reward of +10). To reach any of these 2 states, the agent can either take a short yet dangerous path (going directly right) or take a long yet safe path (going up, then right, then down). The shorter path is dangerous since it extends alongside a row of terminal states with a penalty of -10 each.

The goal of this question is to select values for these 3 parameters (action noise, discount factor and living reward) to control the policy.

- The action noise is the probability that the actual direction of movement ends up being along one of the 2 directions orthogonal to the desired direction. For example, if the noise is 0.2 and the action is up, then the agent has an (1-0.2)=0.8 chance to actually go up and (0.2/2)=0.1 to go left and (0.2/2)=0.1 to go right.
- The discount factor is the gamma as described in the bellman equation.
- The living reward is the reward for going to a non-terminal state. This reward can be positive or negative.

In the file options.py, there are 6 functions question2_1 to question2_6 where each of them returns a dictionary containing the 3 parameters described above. The goal is to select values for there parameters such that the policy behaves as follows:

- 1. For question2_1, we want the policy to seek the near terminal state (reward +1) via the short dangerous path (moving besides the row of -10 state).
- 2. For question2_2, we want the policy to seek the near terminal state (reward +1) via the long safe path (moving away from the row of -10 state).
- 3. For question2_3, we want the policy to seek the far terminal state (reward +10) via the short dangerous path (moving besides the row of -10 state).
- 4. For question2_4, we want the policy to seek the far terminal state (reward +10) via the long safe path (moving away from the row of -10 state).
- 5. For question2_5, we want the policy to avoid any terminal state and keep the episode going on forever.
- 6. For question2_6, we want the policy to seek any terminal state (even ones with the -10 penalty) and try to end the episode in the shortest time possible.

Problem 3: SARSA

Inside reinforcement_learning.py, modify the functions marked by a **TODO** to complete the SARSALearningAgent class.

$$Q(s,a) \leftarrow Q(s,a) + \alpha(r + \gamma Q(s',a') - Q(s,a))$$

Problem 4: Q-Learning

Inside reinforcement_learning.py, modify the functions marked by a **TODO** to complete the QLearningAgent class.

$$Q(s,a) \leftarrow Q(s,a) + \alpha(r + \gamma \max_{a'} Q(s',a') - Q(s,a))$$

Problem 4: Approximate Q-Learning

Inside reinforcement_learning.py, modify the functions marked by a **TODO** to complete the ApproximateQLearningAgent class.

$$w_{ia} \leftarrow w_{ia} + lpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))w_i$$

where w_{ia} is the the weight of the feature x_i in Q(s,a) and $x_1, x_2, ..., x_n$ are the features of the state s. Thus the approximate Q-function can be written as follows:

$$Q(s,a) = \sum_i w_{ia} * x_i$$

Time Limit

In case your computer has a different speed compared to mine (which I will use for testing the submissions), you can use the following information as a reference: On my computer, running speed_test.py takes ~17 seconds. You can measure the run time for this operation on computer by running:

python speed test.py

If your computer is too slow, you can increase the time limit in the testcases files.

Delivery

The delivery deadline is Thursday January 6th 2022 23:59. It should be delivered on **Blackboard**. This is an individual assignment. The delivered code should be solely written by the student who delivered it. Any evidence of plagiarism will lead to receiving **zero** points.