

CS 4/5789 - Programming Assignment 1

February 1, 2023

This assignment is due on **February 15, 2023** at 11:59pm.

Section 0: Requirements

This assignment uses Python 3+. Please install NumPy and Scikit-Learn. For the visualization (described in section 8), also install Pillow, and Matplotlib. If you are on Linux, you will need to install a gui-backend in order for the visualization to work such as PyQt5.

Section 1: Deliverables

For this assignment, we will ask you to work in a 4x4 gridworld. This project consists of the following 5 files

- `MDP.py`: This contains the MDP class which is used to represent our gridworld MDP. This contains the states, actions, rewards, transition probabilities, and the discount factor.
- `visualize.py`: This file can be used to help visualize the V-function and policy. There are also functions that can be used for stepping through value and policy iteration.
- `test.py`: Test case files. Run this with `python test.py`
- `DP.py`: **Student code goes here.** This file will be used for implementing value iteration, exact policy iteration, and approximate policy iteration.
- `IL.py`: **Student code goes here.** This file will be used for all imitation learning methods.

Please implement all the functions with a **TODO** comment in **DP.py** and **IL.py**.

For submission to Gradescope, please do the following:

- Run the visualization for value iteration to iteration 20. Save the figure showing the Q-values. Additionally, save the figure showing the max policy and its corresponding value function.
- Run exact and iterative policy iteration for 5 iterations. Save both figures.
- Run behavior cloning with $N = 75$ datapoints using the code provided to you in `python test.py` and save the figure corresponding to the learned policy.
- Create a PDF writeup in which you include the 5 figures above and also describe the differences between the expert policy and the behavior policy in the IL section.
- Zip the PDF writeup along with the 5 python files listed above into `submit.zip` and submit on Gradescope.

Section 2: Understanding the MDP

2.1 States and rewards

There are 17 (0-16) states where each box corresponds to a state as shown below

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

16

State 15 (the green box) shown above is the “goal” state and has a reward of 200. All the white boxes have reward -1 while the orange boxes are “bad” states that have -80 reward. While our gridworld is 4x4, we need an extra state (state 16) to be the terminal state. This state has reward 0 and is absorbing which means once we enter the terminal state, we cannot leave as any action takes us back to the terminal state. Additionally, once we hit the goal state, we automatically transition to the terminal state. This terminal state allows us to model tasks which terminate at a finite time using the infinite horizon MDP formulation.

Note that the rewards here are outside the range [0,1]. Although many theoretical results assume the reward lies within [0,1], in practice the reward can be any scalar value.

Since we are using a simple state representation, we only need to store the number of states in the MDP class. The rewards, R , is a $|A| \times |S|$ numpy array (note the first index is the action unlike in class). That is, $r(s,a) = R[a,s]$. The rewards in this MDP are deterministic.

2.2 Actions

There are 4 actions that can be taken in each state. These actions are up (0), down (1), left (2), right (3).

2.3 Transition probabilities

The transition probabilities are easy to define for a gridworld MDP which makes it a good setting for implementing value iteration, policy evaluation, and policy iteration. The transition matrix P , in `MDP.PY`, contains the transition probabilities. Note that the transitions are stochastic (i.e. moving right from state 0 isn’t guaranteed to move you to state 1, you may “slip” and move to state 4 or stay in the same spot). In particular, when transitioning from a white or orange box, the agent ends up in the correct spot with probability 0.7, and slips laterally with probability 0.15 to either of the adjacent spots.

P is an $|A| \times |S| \times |S|$ NumPy array. The transition probability $P(s'|s, a)$ can be found by indexing $P[a, s, s']$

2.4 Policies

We will work with deterministic policies and will represent policies using a NumPy array of size $|S|$. The s^{th} entry in such a vector represents the (deterministic) action taken by the policy at state s .

Section 3: Value Iteration with the Q-function

3.1 TODOs

For value iteration, there are 6 functions (TODOs) to complete.

- `computeVfromQ`
- `computeQfromV`
- `extractMaxPifromQ`
- `extractMaxPifromV`
- `valueIterationStep`
- `valueIteration`

3.2 Switching between Q and V functions

Given V^π and Q^π , please implement `computeVfromQ` and `computeQfromV` using the equations

$$\begin{aligned}Q^\pi(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} V^\pi(s'), \\V^\pi(s) &= Q^\pi(s, \pi(s)).\end{aligned}$$

3.3 Computing max policy

From class, we know that if we have the optimal Q-function, Q^* , we can find the optimal policy

$$\pi^*(s) = \arg \max_a Q^*(s)$$

We can also use this operation at the end of value iteration on Q^t (our final Q values) and also during the policy improvement stage of policy iteration. Please implement `extractMaxPifromQ` and `extractMaxPifromV`.

3.4 Value Iteration

First, implement `valueIterationStep`. This implements the following equation

$$\forall s, a : Q^{t+1}(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \max_{a'} Q^t(s', a')$$

Q^t is passed in to this function while the rewards and transition probabilities can be accessed from the attributes. While not strictly necessary, try implementing this function without for-loops.

Next, implement `valueIteration` following the specification. This should consist of running `valueIterationStep` in a loop until the threshold is met. Then, extract the policy from the final Q-function and compute the corresponding V-function. Return the policy, V-function, iterations required until convergence, and the final epsilon.

Section 4: Policy Evaluation

In this section, we ask you to implement both exact policy evaluation and iterative policy evaluation.

4.1 TODOs

For value iteration, there are 2 functions (TODOs) to complete.

- `exactPolicyEvaluation`
- `approxPolicyEvaluation`

4.2 Exact Policy Evaluation

The Bellman equation for deterministic policies, π , from class is as follows

$$\begin{aligned} V^\pi(s) &= r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, \pi(s))} V^\pi(s') \\ &= r(s, \pi(s)) + \gamma \sum_{s'} \mathbb{P}(s'|s, \pi(s)) V^\pi(s') \end{aligned}$$

Using vector notation, this can be written as the linear equation

$$V = R^\pi + \gamma P^\pi V$$

Hint: Use `extractRpi` and `extractPpi` to easily extract $r(s, \pi(s))$ and $P(s'|s, \pi(s))$ in vector form.

By solving this, we can compute V^π . Do **not** use matrix inversion for this. Instead use `np.linalg.solve`.

4.3 Iterative Policy Evaluation

Implement approximate policy evaluation following the algorithm given in class. Each step of iterative policy evaluation involves computing V^{t+1} given the current estimate V^t

$$V^{t+1} = R^\pi + \gamma P^\pi V^t$$

Run this algorithm until convergence (as defined by the tolerance).

Section 5: Policy Iteration

5.1 TODOs

For policy iteration, there are 2 functions (TODOs) to complete.

- `policyIterationStep`
- `policyIteration`

Finally, we will move onto policy iteration. This is another iterative algorithm in which we cycle between policy evaluation and policy improvement to continuously improve our policy. Note the difference between this algorithm and value iteration. In value iteration, we did not compute a policy until we computed the final estimate of the optimal Q-function.

First, implement the function `policyIterationStep`. Given the current policy π^t , compute the value-function for π^t using policy evaluation. This can be done using either exact or iterative policy evaluation. With V^{π^t} , compute the new policy.

Once `policyIterationStep` is implemented, implement `policyIteration` in the same style as `valueIteration`. That is, continuously call `policyIterationStep` until convergence. Convergence in this case occurs when the policy improvement step does not change our current policy.

Section 6: Imitation Learning

In this section, we will ask you to implement an imitation learning algorithm (behavioral cloning) for the gridworld setting. Recall that the goal in imitation learning is, given a set of N expert datapoints $\mathcal{D}^* = \{s_i^*, a_i^*\}_{i=1}^N$ from some unknown, black-box expert policy π^* , to learn a policy π that is as good as the expert.

You will be implementing behavioral cloning (BC), which is the simplest form of imitation learning. In particular, given the dataset, we just perform supervised learning over the dataset in question. More formally, given a set of policies Π , we aim to find the MLE policy

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \frac{1}{N} \sum_{(s^*, a^*) \in \mathcal{D}^*} \log \pi(a^* | s^*).$$

In practice, this is done by minimizing the cross-entropy loss between the actual actions a^* in our expert dataset and the actions predicted by our learner π .

For this assignment, we will be using the policy that you learned in Section 5 using exact policy iteration as the expert policy π^* . We will then collect the expert dataset \mathcal{D}^* and then learn a new behavior policy using behavioral cloning.

6.1 TODOs

For behavioral cloning, there are 2 functions to complete.

- `collectData`
- `trainModel`

6.2 Data Collection

In this section, you will implement a method `collectData` that takes as input the expert policy π^* (which is the policy determined using policy iteration) and the desired number of datapoints N , and returns a training dataset of the form $\mathcal{D}^* = \{s_i^*, a_i^*\}_{i=1}^N$.

Note that in order to train an effective model, we will need to choose a useful feature representation of the input states s . Please store **one-hot** encodings of the input states s in your dataset.

6.3 Training a model

Here you will implement a method `trainModel` that takes the training data from the previous method, and trains a multi-class classification model using behavior cloning.

The policy model class Π that we will be using is the set of classifiers produced by logistic regression. In particular, please use `sklearn.linear_model.LogisticRegression` to train a logistic regression model using supervised learning (with the default cross-entropy loss function).

6.4 Writeup

Run behavior cloning with $N = 75$ datapoints using the code provided to you in `python test.py` and include the generated figure corresponding to the learned policy (and corresponding value function) in your writeup.

Compare the visualization of the expert policy π (which you generated in the previous section for exact PI) with the visualization of the learned BC policy, in particular looking at states other than the starting state 0. What differences do you observe between the expert policy and learned policy, and also between the expert value function and the learned value function?

Section 7: Testing

We have written some basic tests for you to test your implementation. These are separated into tests for VI, PE, PI, and IL. Please note that these tests are basic and do not guarantee that your implementation is correct.

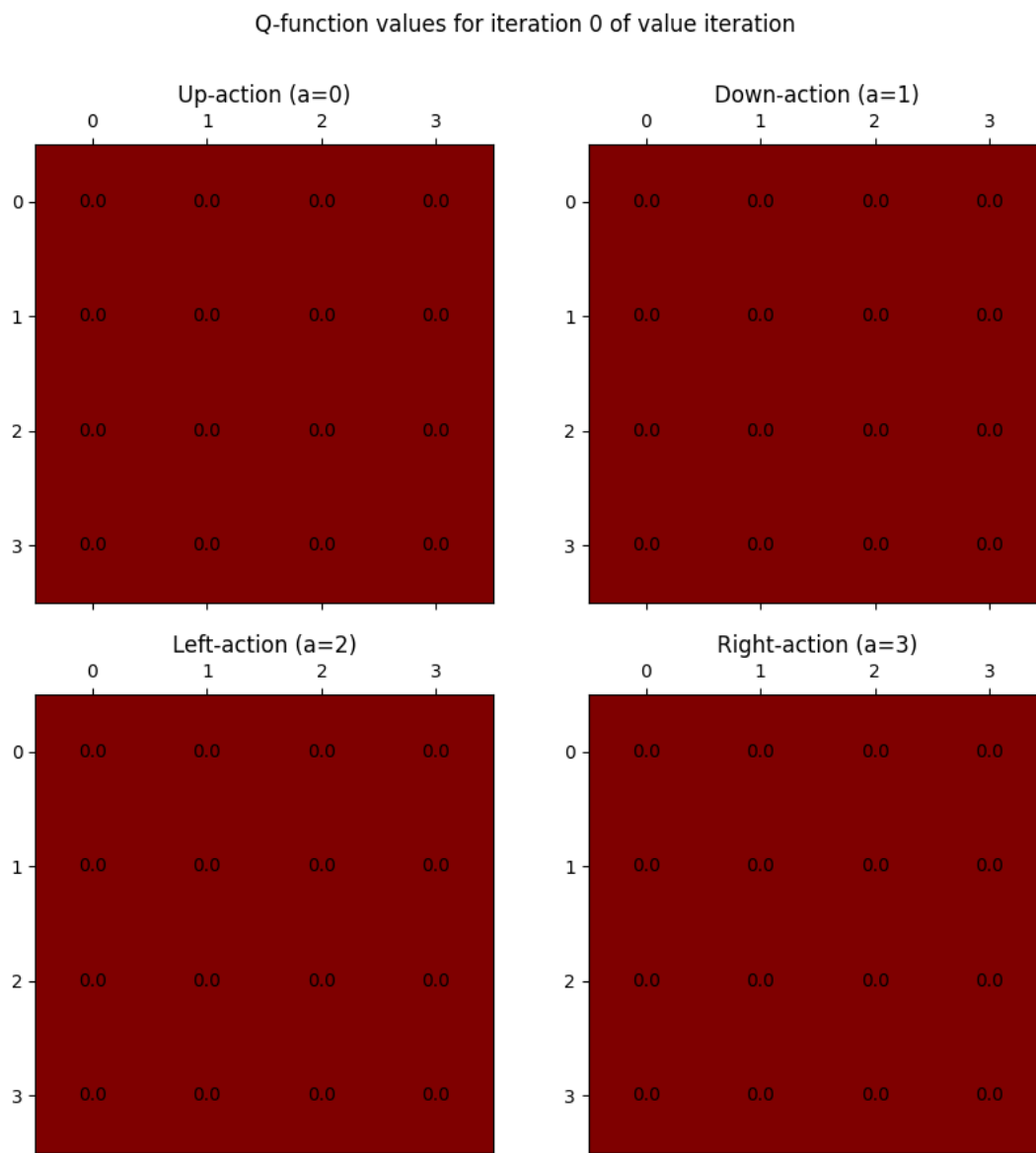
Section 8: Visualization

We have also given you code for visualizing value iteration and policy iteration.

8.1 Value Iteration

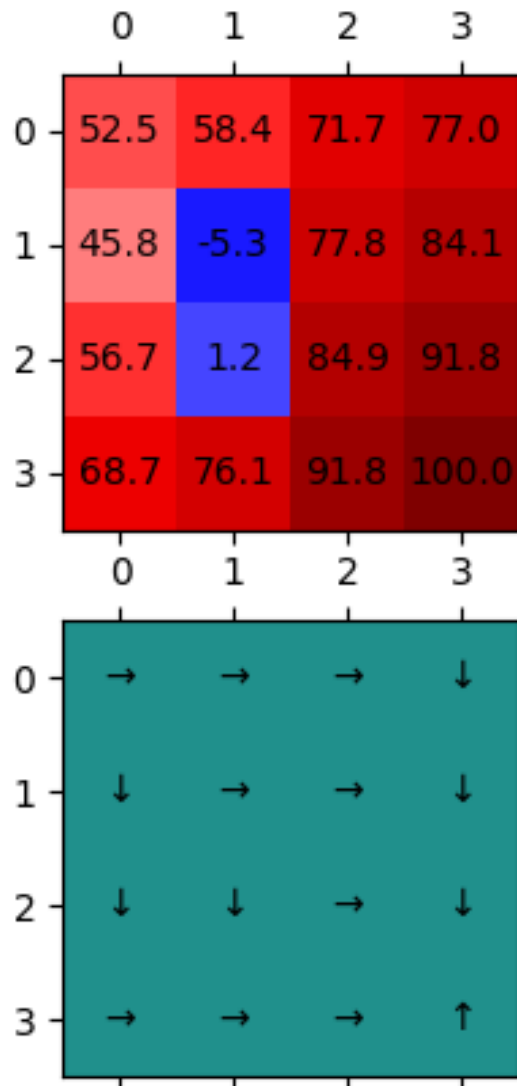
After implementing value iteration, run `python visualize.py VI`. This will allow you to visualize value iteration from an initial Q^0 of all zero's. To change the initial Q^0 , change line 168 in `visualize.py`.

Running this command should open the following window.



This shows the Q-values for each of the 4 actions. To step through an iteration of value iteration, press the **right** arrow key. The grids should update as you step through value iteration.

Normally, value iteration runs until convergence and then the final policy is extracted. However, we have added extra visualization to allow you to see how the policy changes. For any iteration, t , press the **enter** key to compute $\pi(s) = \arg \max_a Q^t(s, a)$ and V^π . A new window should show up showing the V-function and policy. See below.



To continue stepping through value iteration, exit this window and continue pressing the **right** arrow key on the main figure window.

8.2 Policy Iteration

To visualize policy iteration, run `python visualize.py PI_exact` or run `python visualize.py PI_approx`. Again, use the **right** arrow key to step through policy iteration. To change π^0 , change line 171 or 174 of `visualize.py`.