

COMP3667: Reinforcement learning practical 4

Dynamic Programming

robert.lieck@durham.ac.uk

1 Overview

Welcome to the fourth reinforcement learning practical. In this practical, we will dive a bit more into the practical details of dynamic programming methods. In particular, we will:

- Use the OpenAI Gym “Frozen Lake” environment as a basic example.
- Implement *policy evaluation* and understand how this algorithm updates state values.
- Implement *policy improvement* to learn better policies.
- Implement *policy iteration* to learn optimal policies.

For this practical, you will need a basic Python environment with `numpy`, `matplotlib`, and OpenAI `gym` (version 0.20.0).

2 Frozen Lake Environment

Install OpenAI `gym` if you have not done this already. You will need to use version 0.20.0:

[</> copy code](#)

```
1 %%capture
2 !pip install setuptools==65.5.0 "wheel<0.40.0"
3 !apt update
4 !pip install 'gym==0.20.0'
```

Do the relevant imports

[</> copy code](#)

```
1 import gym
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import matplotlib.font_manager
```

and load the Frozen Lake environment

[</> copy code](#)

```
1 name = 'FrozenLake-v1' # small version
2 # name = 'FrozenLake8x8-v1' # larger version
3 env = gym.make(name, is_slippery=False)
4 env.seed(742)
5 env.action_space.seed(742)
6 # named actions
7 LEFT, DOWN, RIGHT, UP = 0, 1, 2, 3
```

The code of the environment is at https://github.com/openai/gym/blob/master/gym/envs/toy_text/frozen_lake.py. Use `slippery=False` for now, because otherwise the optimal solutions are not very intuitive (with `slippery=True`, when the agent tries to move forward it will with equal probability instead slip sideways, but never backwards; therefore it is sometimes better to try to move “away from a hole” instead of “towards the goal”). Use this helper function to visualise the environment in a nice way:

[</> copy code](#)

```
1 def plot(env, v=None, policy=None, col_ramp=1, dpi=175, draw_vals=False, mark_ice=True):
2     # set up plot
3     plt.rcParams['figure.dpi'] = dpi
4     plt.rcParams.update({'axes.edgecolor': (0.32,0.36,0.38)})
```

```

5 plt.rcParams.update({'font.size': 4 if env.env.nrow == 8 else 7})
6 gray = np.array((0.32,0.36,0.38))
7 plt.figure(figsize=(3, 3))
8 ax = plt.gca()
9 ax.set_xticks(np.arange(env.env.ncol)-.5)
10 ax.set_yticks(np.arange(env.env.nrow)-.5)
11 ax.set_xticklabels([])
12 ax.set_yticklabels([])
13 plt.grid(color=(0.42,0.46,0.48), linestyle=':')
14 ax.set_axisbelow(True)
15 ax.tick_params(color=(0.42,0.46,0.48),
16               which='both', top='off', left='off', right='off', bottom='off')
17 # use zero value as dummy if not provided
18 if v is None:
19     v = np.zeros(env.nS)
20 # plot values
21 plt.imshow(1-v.reshape(env.env.nrow,env.env.ncol)**col_ramp,
22           cmap='gray', interpolation='none',
23           clim=(0,1), zorder=-1)
24 # go through states
25 for s in range(env.nS):
26     x, y = s % env.env.nrow, s // env.env.ncol
27     # print numeric values
28     if draw_vals and v[s] > 0:
29         vstr = '{0:.1e}'.format(v[s]) if env.env.nrow == 8 else '{0:.6f}'.format(v[s])
30         plt.text(x - 0.45, y + 0.45, vstr, color=(0.2, 0.8, 0.2), fontname='Sans')
31     # mark ice, start, goal
32     if env.desc.tolist()[y][x] == b'F':
33         plt.text(x-0.45,y-0.3, 'ice', color=(0.5, 0.6, 1), fontname='Sans')
34         if mark_ice:
35             ax.add_patch(plt.Circle((x, y), 0.2, color=(0.7, 0.8, 1), zorder=0))
36     elif env.desc.tolist()[y][x] == b'S':
37         plt.text(x-0.45,y-0.3, 'start',color=(0.2,0.5,0.5), fontname='Sans',
38               weight='bold')
39     elif env.desc.tolist()[y][x] == b'G':
40         plt.text(x-0.45,y-0.3, 'goal', color=(0.7,0.2,0.2), fontname='Sans',
41               weight='bold')
42         continue # don't plot policy for goal state
43     else:
44         continue # don't plot policy for holes
45     # plot policy
46     def plot_arrow(x, y, dx, dy, v, scale=0.4):
47         plt.arrow(x, y, scale * float(dx), scale * float(dy), color=gray+0.2*(1-v),
48               head_width=0.1, head_length=0.1, zorder=1)
49     if policy is not None:
50         a = policy[s]
51         if a[0] > 0.0: plot_arrow(x, y, -a[0], 0., v[s]) # left
52         if a[1] > 0.0: plot_arrow(x, y, 0., a[1], v[s]) # down
53         if a[2] > 0.0: plot_arrow(x, y, a[2], 0., v[s]) # right
54         if a[3] > 0.0: plot_arrow(x, y, 0., -a[3], v[s]) # up
55 plt.show()

```

Have a look

[</> copy code](#)

```

1 print('action space: ' + str(env.action_space))
2 print('reward range: ' + str(env.reward_range))
3 print('observation space: ' + str(env.observation_space))
4 plot(env=env)
5 --> action space: Discrete(4)
6 --> reward range: (0, 1)
7 --> observation space: Discrete(16)

```

- Define a uniform policy and plot it in the environment. *Hint:* `env.nS` and `env.nA` give you the number of states and actions, respectively; you can provide the policy to the plotting function via the `policy` argument.

3 Policy Evaluation

Now, we would like to know how well a policy performs, that is, what the state value of a particular state is when following the policy in the future. This *policy evaluation* procedure can be efficiently done using dynamic programming, which iteratively improves state value estimates and converges to the true state values.

- Write a `policy_eval_step` function that takes an initial state value estimate $v_\pi^k(s)$ and computes an improved estimate $v_\pi^{k+1}(s)$ as

$$v_\pi^{k+1}(s) = \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) [\mathcal{R}_{ss'}^a + \gamma v_\pi^k(s')] . \quad (1)$$

Hint: `env.env.P[s][a]` is giving you a list of tuples `(p, s', r, done)`, one for each possible transition from state `s` when taking action `a`: `p` is the probability of this transition to happen, `s'` is the state the agent transitions to, `r` is the reward it receives, and `done` indicates whether the episode is over (which you will not need to use).

- Initialise the state values with zero, take one `policy_eval_step` at a time and plot the result to observe how updates are being performed. *Hint:* You can provide the state values to the plot function using the `v` argument, specifying `draw_vals=True` additionally shows the numeric state values.
- On an intuitive level, how would you describe the dynamics you see? What seems to be inefficient about the current implementation? How could this be improved?
- Implement a modified `policy_eval_step_inplace` version of the function that performs state value updates in-place. That is, instead of clearly separating the “old” values $v_\pi^k(s)$ and the “new” values $v_\pi^{k+1}(s)$, you operate on a single state value estimate $v_\pi(s)$, which is updated as you go. *Hint:* Make sure the updates for a particular state are still “atomic” and you do not use the half-computed values (in case of transitions that *stay* in a particular state).
 - Think about a clever order in which to update state values in-place. *Hint:* States are ordered from top-left (start: 0) to bottom-right (goal: 15).
 - Again, observe the step-wise update of values from one iteration to the next. How does that compare to the original implementation?
- Write a `policy_evaluation` function that iteratively updates the state values using the `policy_eval_step` or `policy_eval_step_inplace` function and stops if they do not change (by some small tolerance value). Print the number of iterations needed to converge and compare for the `policy_eval_step` and `policy_eval_step_inplace` implementation.
- How many iterations do you need until state values have converged to their true value? To make it simpler, do the following though experiment: Take an environment that has only a single state and a single action (so nothing can really change and there is only one possible policy) and you get a reward of 1 upon every transition. Take the update equation for the state value from above, which now simplifies to

$$v^{k+1} = 1 + \gamma v^k . \quad (2)$$

Can you write down v^n in a non-recursive form assuming you start with $v^0 = 0$? Can you write down v^∞ in closed form? Is v^∞ the exact state value? How long does it take to converge? What happens for $\gamma = 1$ as opposed to $\gamma < 1$?

4 Policy Improvement

- Implement a function that computes state-action values $q_\pi(s, a)$ (for all actions in a given state) from the state values $v_\pi(s)$ using their known relation

$$q_\pi(s, a) = \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma v_\pi(s')] . \quad (3)$$

- Implement a `policy_improvement` function that takes a state value estimate and defines a policy (by first computing state-action values) that achieves maximal return (i.e. always chooses an action with maximal state-action value). Optionally, either choose actions deterministically, or choose all actions with maximum value with the same probability (if there is more than one such action).
- Load the larger 'FrozenLake8x8-v1' environment (again with `is_slippery=False` for now), compute state values by evaluating the uniform policy, plot the result. Then compute an improved policy and plot the result again.

5 Policy Iteration

- Starting from the improved policy from above, perform two updates by doing
 - policy evaluation (plot the results)
 - policy improvement (plot the results).

Use a stochastic policy improvement (i.e. choose optimal action with equal probability) and a discount value of $\gamma = 1$. What do you observe? Do you see anything that could be problematic? Can you explain what you observe (you may need to print the raw state values)? Would a deterministic policy improvement help? *Hint:* Remember what we learned above about value convergence and paths of finite/infinite length.

- Change the value to $\gamma = 0.999$ and do another two updates. What is different? *Hint:* You can use $\gamma = 0.9$ to see the effect more clearly.
