





```
[INFO] In 800 iteration, current mean episode reward is 0.867.
[INFO] In 900 iteration, current mean episode reward is 0.867.
[INFO] In 1000 iteration, current mean episode reward is 0.867.
[INFO] In 1100 iteration, current mean episode reward is 0.867.
[INFO] In 1200 iteration, current mean episode reward is 0.867.
[INFO] In 1300 iteration, current mean episode reward is 0.867.
[INFO] In 1400 iteration, current mean episode reward is 0.867.
[INFO] In 1500 iteration, current mean episode reward is 0.867.
[INFO] In 1600 iteration, current mean episode reward is 0.867.
[INFO] In 1700 iteration, current mean episode reward is 0.867.
[INFO] In 1800 iteration, current mean episode reward is 0.867.
We found policy is not changed anymore at iteration 1800. Current mean episode reward
is 0.867. Stop training.
```

```
In [17]: # Run this cell without modification

print("Your value iteration agent achieve {} mean episode reward. The optimal score "
      "should be almost {}.".format(vi_agent.evaluate(), 0.86))
```

Your value iteration agent achieve 0.867 mean episode reward. The optimal score should be almost 0.86.

```
In [18]: # Run this cell without modification

vi_agent.render()

(Right)
FFFFFFF
FFFFFFF
FFFFFFF
FFFFFFF
FFFFFFF
FFFFFFF
```

```
In [19]: # Run this cell without modification

vi_agent.print_table()

+-----+-----+-----+-----+-----+-----+-----+-----+
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
+-----+-----+-----+-----+-----+-----+-----+
| 0 | 1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 1 | 1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 2 | 1.000|0.978|0.926|0.000|0.857|0.946|0.982|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 3 | 1.000|0.935|0.801|0.475|0.624|0.000|0.945|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 4 | 1.000|0.826|0.542|0.000|0.539|0.611|0.852|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 5 | 1.000|0.000|0.000|0.168|0.383|0.442|0.000|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 6 | 1.000|0.000|0.195|0.121|0.000|0.332|0.000|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 7 | 1.000|0.732|0.463|0.000|0.277|0.555|0.777|0.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
```

Congratulation! You have successfully implemented the value iteration trainer (if and only if no error happens at the above cells). Few further problems for you to investigate:

- Do you see that some iteration during training yields better rewards than the final one? Why does that happen?
- What is the impact of the discount factor gamma?
- What is the impact of the value function convergence criterion epsilon?

If you are interested in doing more investigation (not limited to these two), feel free to open new cells at the end of this notebook and left a clear trace of your thinking and coding, which leads to extra credit if you do a good job. It's an optional job, and you can ignore it.

Now let's continue our journey!

### Section 2.3: Compare two model-based agents

Now we have two agents: `pi_agent` and `vi_agent`. They are believed to be the optimal policy in this environment. Can you compare the policy of two of them and use a clean and clear description or figures to show your conclusion?

# Solve the TODO and remove `pass # TODO` try to compare two trained agents' policies # hint: `trainer.print_table()` may give you a better sense. The two cells below shows the final value table for each of the 64 states derived from policy iteration and value iteration respectively. Every corresponding state of the two tables are the same. As for each of the policies of the two agents, it would choose the action that leads to a next state that yields the highest value. Therefore, the policies of the two agents have to be identical.

```
In [20]: pi_agent.print_table()

+-----+-----+-----+-----+-----+-----+-----+-----+
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
+-----+-----+-----+-----+-----+-----+-----+
| 0 | 1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 1 | 1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 2 | 1.000|0.978|0.926|0.000|0.857|0.946|0.982|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 3 | 1.000|0.935|0.801|0.475|0.624|0.000|0.945|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 4 | 1.000|0.826|0.542|0.000|0.539|0.611|0.852|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 5 | 1.000|0.000|0.000|0.168|0.383|0.442|0.000|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 6 | 1.000|0.000|0.195|0.121|0.000|0.332|0.000|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 7 | 1.000|0.732|0.463|0.000|0.277|0.555|0.777|0.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
```

```
In [21]: vi_agent.print_table()

+-----+-----+-----+-----+-----+-----+-----+-----+
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
+-----+-----+-----+-----+-----+-----+-----+
| 0 | 1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 1 | 1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 2 | 1.000|0.978|0.926|0.000|0.857|0.946|0.982|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 3 | 1.000|0.935|0.801|0.475|0.624|0.000|0.945|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 4 | 1.000|0.826|0.542|0.000|0.539|0.611|0.852|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 5 | 1.000|0.000|0.000|0.168|0.383|0.442|0.000|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 6 | 1.000|0.000|0.195|0.121|0.000|0.332|0.000|1.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
| 7 | 1.000|0.732|0.463|0.000|0.277|0.555|0.777|0.000|
| | | | | | | | | |
+-----+-----+-----+-----+-----+-----+-----+
```

# You can do more investigation here if you wish. Leave it blank if you don't. "" hint: generating random action at each call of policy may lead to a failure of convergence, try generate random actions at initialization # and fix it during the training. Why is this the case? ""

### Conclusion and Discussion

In this assignment, we learn how to use Gym package, how to use Object Oriented Programming idea to build a basic tabular RL algorithm.

It's OK to leave the following cells empty. In the next markdown cell, you can write whatever you like. Like the suggestion on the course, the confusing problems in the assignments, and so on.

If you want to do more investigation, feel free to open new cells via `Esc + B` after the next cells and write codes in it, so that you can reuse some result in this notebook. Remember to write sufficient comments and documents to let others know what you are doing.

Following the submission instruction in the assignment to submit your assignment to our staff. Thank you!

```
...
```

```
In [ ] :
```