

```
Last login: Sat Jun  6 14:34:29 on ttys001
cd /Users/liuxuemin/Desktop/NDResearch/Machine\ Learning\ courses\ outside\ ND/S
tanford-CS234-ReinforcementLearning/assignment1; clear; pwd
(base) Xuemins-MacBook-Pro:~ liuxuemin$ cd /Users/liuxuemin/Desktop/NDResearch/M
achine\ Learning\ courses\ outside\ ND/Stanford-CS234-ReinforcementLearning/assi
gnment1; clear; pwd
```

```
/Users/liuxuemin/Desktop/NDResearch/Machine Learning courses outside ND/Stanford
-CS234-ReinforcementLearning/assignment1
(base) Xuemins-MacBook-Pro:assignment1 liuxuemin$ python vi_and_pi.py
```

Results of the Deterministic environment

Beginning Policy Iteration

```
SFFF
FHFH
FFFH
HFFG
  (Down)
SFFF
FHFH
FFFH
HFFG
  (Down)
SFFF
FHFH
FFFH
HFFG
  (Right)
SFFF
FHFH
FFFH
HFFG
  (Down)
```

SFFF
FHFH
FFFH
HFFG
 (Right)

SFFF
FHFH
FFFH
HFFG
 (Right)

SFFF
FHFH
FFFH
HFFG

Episode reward: 1.000000

Beginning Value Iteration

SFFF
FHFH
FFFH
HFFG
 (Down)

SFFF
FHFH
FFFH
HFFG
 (Down)

SFFF
FHFH
FFFH
HFFG
 (Right)

SFFF
FHFH
FFFH
HFFG
 (Down)

SFFF
FHFH
FFFH
HFFG
 (Right)

SFFF
FHFH
FFFH
HFFG
 (Right)

SFFF
FHFH

FFFH

HFFG

Episode reward: 1.000000

(base) Xuemins-MacBook-Pro:assignment1 liuxuemin\$ python vi_and_pi.py

Results of the Stochastic environment

Beginning Policy Iteration

SFFF

FHFH

FFFH

HFFG

(Left)

SFFF

FHFH

FFFH

HFFG

(Left)

SFFF

FHFH

FFFH

HFFG

(Left)

SFFF

FHFH

FFFH

HFFG

(Left)

SFFF

FHFH

FFFH

HFFG

(Left)

SFFF

FHFH

FFFH

HFFG

(Left)

SFFF

FHFH

FFFH

HFFG

(Left)

SFFF

FHFH

FFFH











HFFG

(Left)

SFFF

FHFH

FFFH

HFFG
 (Left)
SFFF
FHFH
FFH
HFFG
 (Up)
SFFF
FHFH
FFH
HFFG
 (Up)
SFFF
FHFH
FFH
HFFG
 (Up)
SFFF
HFH
FFFH
HFFG
 (Left)
SFFF
HFH
FFFH
HFFG
 (Left)
SFFF
FHFH
FFH
HFFG
 (Up)
SFFF
FHFH
FFH
HFFG
 (Up)
SFFF
FHFH
FFH
HFFG
 (Up)
SFFF
HFH
FFFH
HFFG
 (Left)
SFFF
HFH
FFFH
HFFG
 (Left)

SFFF

FHFH

FFH

HFFG

(Up)

SFFF

HFH

FFFH

HFFG

(Left)

SFFF

FHFH

FFH

HFFG

(Up)

SFFF

FHFH

FFH

HFFG

(Down)

SFFF

FHFH

FFFH

FFG

(Right)

SFFF

FHFH

FFFH

FG

(Down)

SFFF

FHFH

FFFH

FG

(Right)

SFFF

FHFH

FFFH

FG

(Down)

SFFF

FHFH

FFFH

FG

(Down)

SFFF

FHFH

FFFH

FG

(Down)

SFFF

FHFH

FFFH
HFFG
Episode reward: 1.000000

Beginning Value Iteration

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG

(Left)
SFFF
FHFH
FFFH
HFFG

(Left)
SFFF
FHFH
FFFH
HFFG

(Left)
SFFF
FHFH
FFFH
HFFG

(Left)
SFFF
FHFH
FFFH
HFFG

(Left)
SFFF
FHFH
FFFH
HFFG

(Left)
SFFF
FHFH
FFFH
HFFG

(Left)
SFFF
FHFH
FFFH
HFFG

(Left)
SFFF
FHFH
FFFH
HFFG

(Left)
SFFF
FHFH
FFFH
HFFG

(Left)
SFFF
FHFH
FFFH
HFFG

(Left)
SFFF

[illegible]

HFFG
 (Left)
SFFF
FHFH
FFFFH
HFFG
 (Left)
SFFF
FHFH
FFFFH
HFFG
 (Left)
SFFF
FHFH
FFFFH
HFFG
 (Left)
SFFF
FHFH
FFFFH
HFFG
 (Left)
SFFF
FHFH
FFFFH
HFFG
 (Left)
SFFF
FHFH
FFFFH
HFFG
 (Left)
SFFF
FHFH
FFFFH
HFFG
 (Left)
SFFF
FHFH
FFFFH
HFFG
 (Left)
SFFF
FHFH
FFFFH
HFFG
 (Left)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Up)

SFFF
FHFH
FFFH
HFFG
(Down)

SFFF
FHFH
FFFH
HFFG
(Up)

SFFF
FHFH
FFFH
HFFG
(Left)

SFFF
FHFH
FFFH
HFFG
(Up)

SFFF
FHFH

FFFH
HFFG
(Up)
SFFF
FHFH
FFFH
HFFG
(Down)
SFFF
FHFH
FFFH
HFFG
(Up)
SFFF
FHFH
FFFH
HFFG
(Down)
SFFF
FHFH
FFFH
HFFG
(Up)
SFFF
FHFH
FFFH
HFFG
(Up)
SFFF
FHFH
FFFH
HFFG
(Down)
SFFF
FHFH
FFFH
HFFG
(Left)
SFFF
FHFH
FFFH
HFFG
(Down)
SFFF
FHFH
FFFH
HFFG
(Right)
SFFF
FHFH
FFFH
HFFG

```
(Down)
SFFF
FHFH
FFFH
HF G
(Down)
SFFF
FHFH
FFFH
HF G
(Down)
SFFF
FHFH
FFFH
HFF G
Episode reward: 1.000000
(base) Xuemins-MacBook-Pro:assignment1 liuxuemin$
```

If you actually run the code, in the case of Deterministic, both Policy Iteration and Value Iteration arrive well as Goals at once.

However, in the stochastic environment, it seems very hesitant and tends to proceed for a long time.

The number of Iterations increased both in the case of Value Iteration and Policy Iteration. This is reasonable.

In addition, it was confirmed that when the environment is stochastic, the Deterministic environment and Optimal Policy are also changed.

If the V and P converge with the same policy, the learning seems to be good.