

# E1

Tae Hyeon Kwon, undergrad(ITM)

2021-01-22

*page 21*

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

R = np.hstack((np.repeat(-1.5,4,axis=0),-0.5,np.repeat(-1.5,2,axis=0),0)).reshape(-1,1).T
states = np.arange(0,80,10)

P=np.matrix([[.1,0,.9,0,0,0,0,0],
             [.1,0,0,.9,0,0,0,0],
             [0,.1,0,0,.9,0,0,0],
             [0,0,.1,0,0,.9,0,0],
             [0,0,0,.1,0,0,.9,0],
             [0,0,0,0,.1,0,0,.9],
             [0,0,0,0,0,.1,0,.9],
             [0,0,0,0,0,0,0,1]])

P = pd.DataFrame(P,columns=states)

print(R)

## [[-1.5 -1.5 -1.5 -1.5 -0.5 -1.5 -1.5  0.  ]]

print(P)

##      0      10      20      30      40      50      60      70
## 0  0.1  0.0  0.9  0.0  0.0  0.0  0.0  0.0
## 1  0.1  0.0  0.0  0.9  0.0  0.0  0.0  0.0
## 2  0.0  0.1  0.0  0.0  0.9  0.0  0.0  0.0
## 3  0.0  0.0  0.1  0.0  0.0  0.9  0.0  0.0
## 4  0.0  0.0  0.0  0.1  0.0  0.0  0.9  0.0
## 5  0.0  0.0  0.0  0.0  0.1  0.0  0.0  0.9
## 6  0.0  0.0  0.0  0.0  0.0  0.1  0.0  0.9
## 7  0.0  0.0  0.0  0.0  0.0  0.0  0.0  1.0
```

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

R = np.hstack((np.repeat(-1.5,4,axis=0),-0.5,np.repeat(-1.5,2,axis=0),0)).reshape(-1,1)
states = np.arange(0,80,10)

P = np.matrix([[.1,0,.9,0,0,0,0,0],
               [.1,0,0,.9,0,0,0,0],
               [0,.1,0,0,.9,0,0,0],
               [0,0,.1,0,0,.9,0,0],
               [0,0,0,.1,0,0,.9,0],
               [0,0,0,0,.1,0,0,.9],
               [0,0,0,0,0,.1,0,.9],
               [0,0,0,0,0,0,0,1]])

gamma = 1.0
epsilon = 10**(-8)
v_old = np.zeros((8,1))
v_new = R + np.dot(gamma*P,v_old)

results = v_old.T
results = np.append(results,v_new.T,axis=0)

while np.max(abs(v_new-v_old))>epsilon:
    v_old = v_new
    v_new = R + np.dot(gamma*P,v_old)
    results = np.append(results,v_new.T,axis=0)

results = pd.DataFrame(results,columns=states)

print(v_new.T)

```

```

## [[-5.80592905 -5.2087811 -4.13926239 -3.47576467 -2.35376031 -1.73537603
##    -1.6735376  0.          ]]

```

```
print(results.head(n=7))
```

```

##          0          10          20          30          40          50          60       70
## 0  0.00000  0.000000  0.000000  0.000000  0.000000  0.00000  0.00000  0.0
## 1 -1.50000 -1.500000 -1.50000 -1.500000 -0.500000 -1.5000 -1.50000  0.0
## 2 -3.00000 -3.000000 -2.10000 -3.000000 -2.000000 -1.5500 -1.65000  0.0
## 3 -3.69000 -4.500000 -3.60000 -3.105000 -2.285000 -1.7000 -1.65500  0.0
## 4 -5.10900 -4.663500 -4.00650 -3.390000 -2.300000 -1.7285 -1.67000  0.0
## 5 -5.61675 -5.061900 -4.03635 -3.456300 -2.342000 -1.7300 -1.67285  0.0
## 6 -5.69439 -5.172345 -4.11399 -3.460635 -2.351195 -1.7342 -1.67300  0.0

```

```
print(results.tail(n=7))
```

```

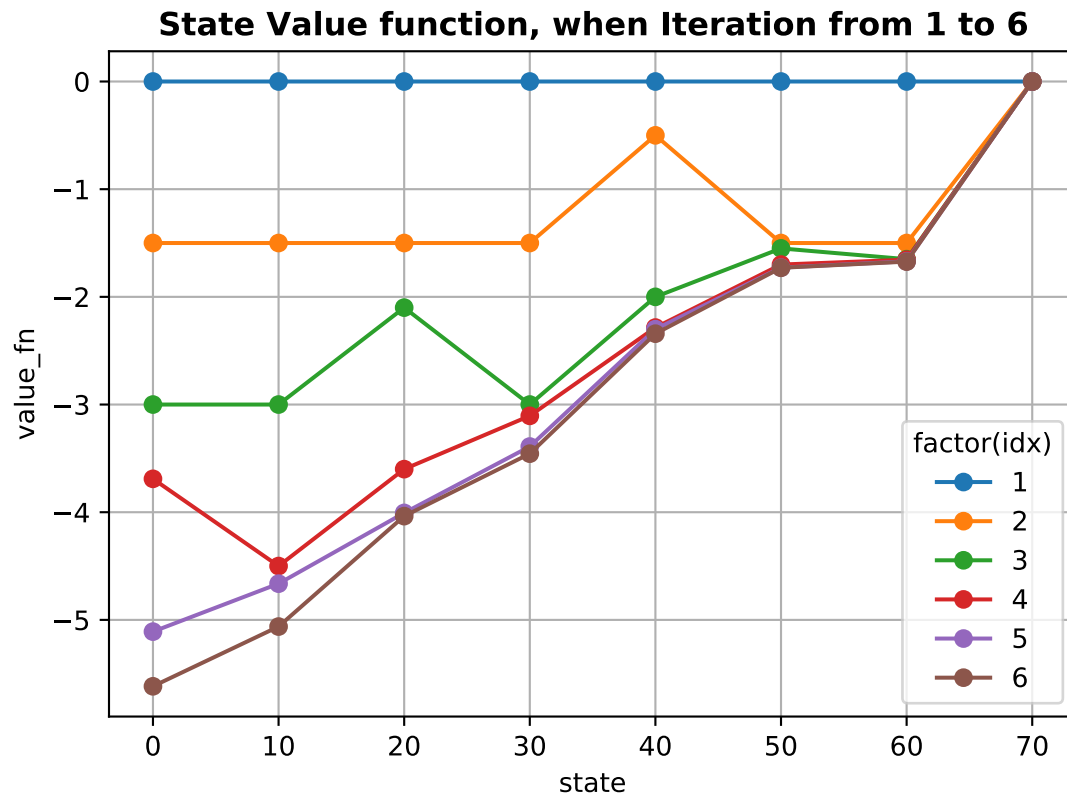
##          0          10          20          30          40          50          60       70

```

```
## 16 -5.805927 -5.208779 -4.139262 -3.475764 -2.35376 -1.735376 -1.673538 0.0
## 17 -5.805928 -5.208781 -4.139262 -3.475765 -2.35376 -1.735376 -1.673538 0.0
## 18 -5.805929 -5.208781 -4.139262 -3.475765 -2.35376 -1.735376 -1.673538 0.0
## 19 -5.805929 -5.208781 -4.139262 -3.475765 -2.35376 -1.735376 -1.673538 0.0
## 20 -5.805929 -5.208781 -4.139262 -3.475765 -2.35376 -1.735376 -1.673538 0.0
## 21 -5.805929 -5.208781 -4.139262 -3.475765 -2.35376 -1.735376 -1.673538 0.0
## 22 -5.805929 -5.208781 -4.139262 -3.475765 -2.35376 -1.735376 -1.673538 0.0
```

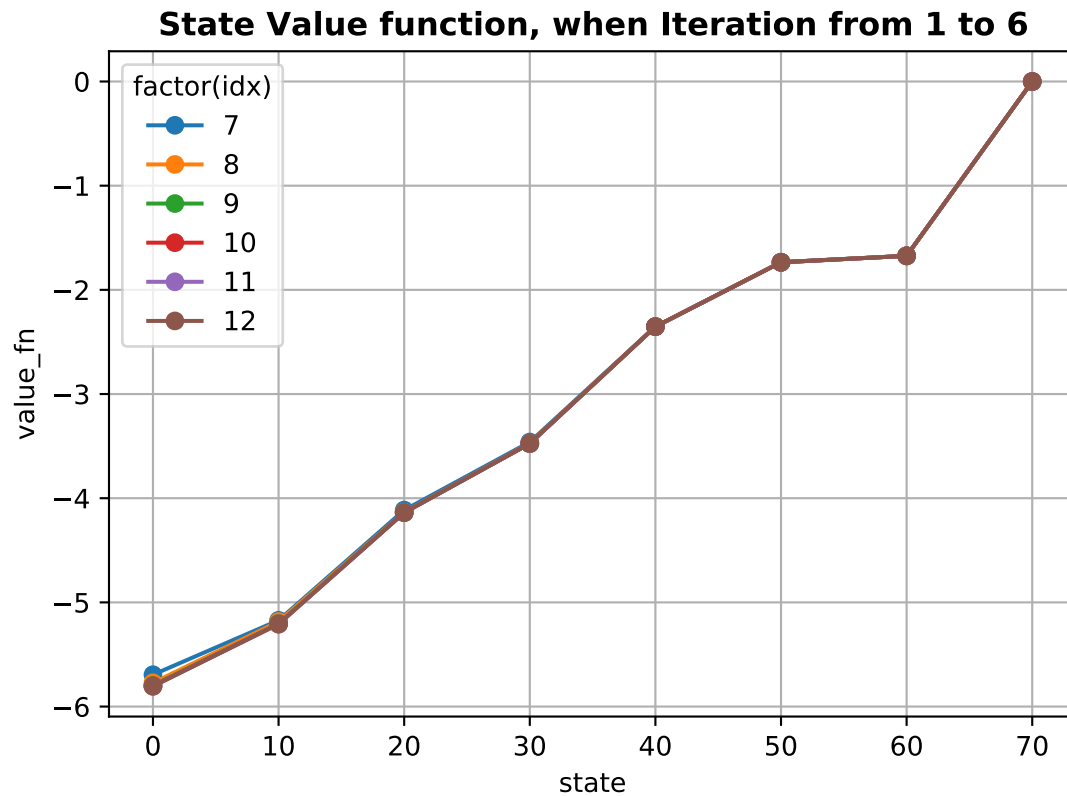
*Iteration from 1 to 6*

```
plt.plot(states,results.iloc[0],marker='o',label='1')
plt.plot(states,results.iloc[1],marker='o',label='2')
plt.plot(states,results.iloc[2],marker='o',label='3')
plt.plot(states,results.iloc[3],marker='o',label='4')
plt.plot(states,results.iloc[4],marker='o',label='5')
plt.plot(states,results.iloc[5],marker='o',label='6')
plt.grid(True)
plt.legend(title='factor(idx)')
plt.xlabel('state')
plt.ylabel('value_fn')
plt.title('State Value function, when Iteration from 1 to 6',fontweight='bold')
plt.show()
```



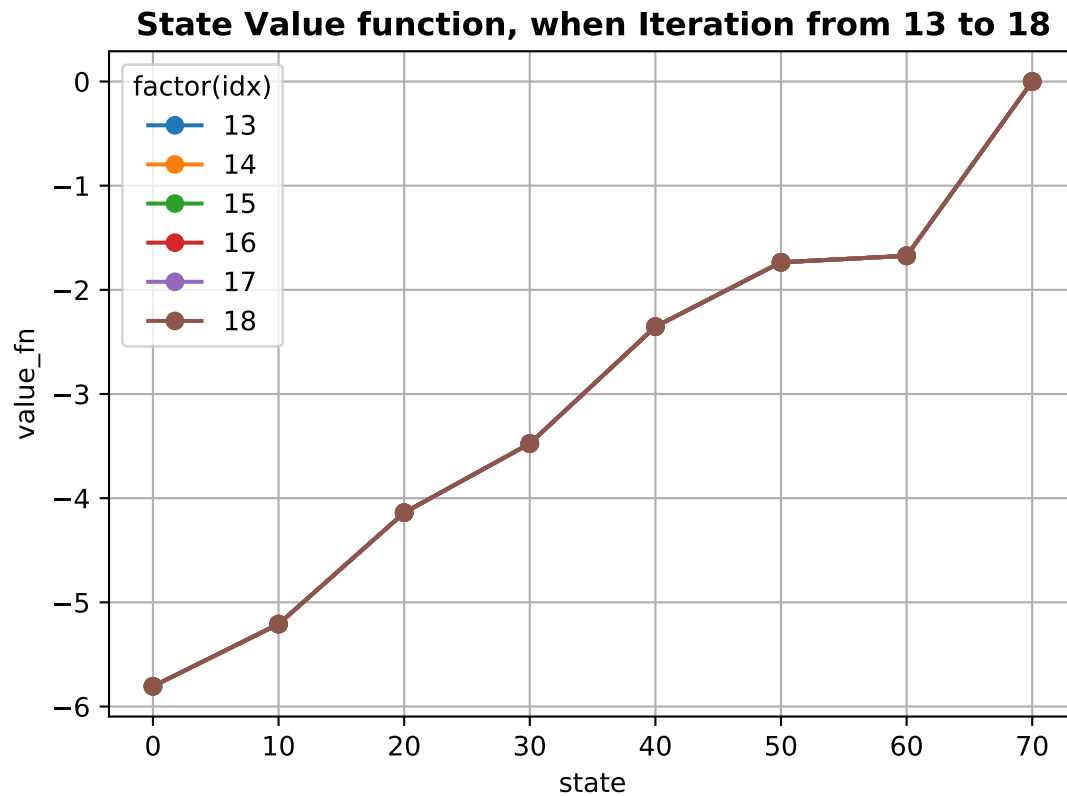
*Iteration from 7 to 12*

```
plt.plot(states,results.iloc[6],marker='o',label='7')
plt.plot(states,results.iloc[7],marker='o',label='8')
plt.plot(states,results.iloc[8],marker='o',label='9')
plt.plot(states,results.iloc[9],marker='o',label='10')
plt.plot(states,results.iloc[10],marker='o',label='11')
plt.plot(states,results.iloc[11],marker='o',label='12')
plt.grid(True)
plt.legend(title='factor(idx)')
plt.xlabel('state')
plt.ylabel('value_fn')
plt.title('State Value function, when Iteration from 1 to 6',fontweight='bold')
plt.show()
```



*Iteration from 13 to 18*

```
plt.plot(states,results.iloc[12],marker='o',label='13')
plt.plot(states,results.iloc[13],marker='o',label='14')
plt.plot(states,results.iloc[14],marker='o',label='15')
plt.plot(states,results.iloc[15],marker='o',label='16')
plt.plot(states,results.iloc[16],marker='o',label='17')
plt.plot(states,results.iloc[17],marker='o',label='18')
plt.grid(True)
plt.legend(title='factor(idx)')
plt.xlabel('state')
plt.ylabel('value_fn')
plt.title('State Value function, when Iteration from 13 to 18',fontweight='bold')
plt.show()
```



## Policy evaluation 2

```
states=np.arange(0,70+10,10).astype('str')

pi_speed=np.c_[np.repeat(0,len(states)),np.repeat(1,len(states))]
pi_speed=pd.DataFrame(data=pi_speed, index=states, columns=['normal','speed'])

print(pi_speed)
```

```
##      normal  speed
## 0         0      1
## 10        0      1
## 20        0      1
## 30        0      1
## 40        0      1
## 50        0      1
## 60        0      1
## 70        0      1
```

```
R_s_a=pd.DataFrame(np.matrix([-1,-1,-1,-1,0.0,-1,-1,0,-1.5,-1.5,-1.5,-1.5,-0.5,-1.5,-1.5,0]).reshape((len(states),2))))
print(R_s_a)
```

```
##      normal  speed
## 0      -1.0  -1.5
## 10     -1.0  -1.5
## 20     -1.0  -1.5
## 30     -1.0  -1.5
## 40      0.0  -0.5
## 50     -1.0  -1.5
## 60     -1.0  -1.5
## 70      0.0   0.0
```

```
def reward_fn(given_pi):
    R_s_a=pd.DataFrame(np.matrix([-1,-1,-1,-1,0.0,-1,-1,0,-1.5,-1.5,-1.5,-1.5,-0.5,-1.5,-1.5,0]).reshape((len(states),2))))

    R_pi=np.asarray((given_pi*R_s_a).sum(axis=1)).reshape((-1,1))

    return R_pi

print(reward_fn(pi_speed))
```

```
## [[-1.5]
##  [-1.5]
##  [-1.5]
##  [-1.5]
##  [-0.5]
##  [-1.5]
##  [-1.5]
##  [ 0. ]]
```

```

P_normal=pd.DataFrame(np.matrix([[0,1,0,0,0,0,0,0],
                                [0,0,1,0,0,0,0,0],
                                [0,0,0,1,0,0,0,0],
                                [0,0,0,0,1,0,0,0],
                                [0,0,0,0,0,1,0,0],
                                [0,0,0,0,0,0,1,0],
                                [0,0,0,0,0,0,0,1],
                                [0,0,0,0,0,0,0,1]]), index=states, columns=states)

P_speed=pd.DataFrame(np.matrix([[.1,0,.9,0,0,0,0,0],
                                [.1,0,0,.9,0,0,0,0],
                                [0,.1,0,0,.9,0,0,0],
                                [0,0,.1,0,0,.9,0,0],
                                [0,0,0,.1,0,0,.9,0],
                                [0,0,0,0,.1,0,0,.9],
                                [0,0,0,0,0,.1,0,.9],
                                [0,0,0,0,0,0,0,1]]), index=states, columns=states)

def transition(given_pi, states, P_normal, P_speed):
    P_out=pd.DataFrame(np.zeros((len(states),len(states))),index=states, columns=states)

    for s in states:
        action_dist=given_pi.loc[s]
        P=action_dist['normal']*P_normal+action_dist['speed']*P_speed
        P_out.loc[s]=P.loc[s]

    return P_out

```



```
print(pi_speed)
```

```
##      normal  speed
## 0         0      1
## 10        0      1
## 20        0      1
## 30        0      1
## 40        0      1
## 50        0      1
## 60        0      1
## 70        0      1
```

```
print(transition(pi_speed, states=states, P_normal=P_normal, P_speed=P_speed))
```

```
##      0    10    20    30    40    50    60    70
## 0  0.1  0.0  0.9  0.0  0.0  0.0  0.0  0.0
## 10 0.1  0.0  0.0  0.9  0.0  0.0  0.0  0.0
## 20 0.0  0.1  0.0  0.0  0.9  0.0  0.0  0.0
## 30 0.0  0.0  0.1  0.0  0.0  0.9  0.0  0.0
## 40 0.0  0.0  0.0  0.1  0.0  0.0  0.9  0.0
## 50 0.0  0.0  0.0  0.0  0.1  0.0  0.0  0.9
## 60 0.0  0.0  0.0  0.0  0.0  0.1  0.0  0.9
## 70 0.0  0.0  0.0  0.0  0.0  0.0  0.0  1.0
```

```
pi_50=pd.DataFrame(np.c_[np.repeat(0.5,len(states)),np.repeat(0.5,len(states))], index=states, columns=states)
print(pi_50)
```

```
##      normal  speed
## 0         0.5    0.5
## 10        0.5    0.5
## 20        0.5    0.5
## 30        0.5    0.5
## 40        0.5    0.5
## 50        0.5    0.5
## 60        0.5    0.5
## 70        0.5    0.5
```

```
print(transition(pi_50, states=states, P_normal=P_normal, P_speed=P_speed))
```

```
##      0    10    20    30    40    50    60    70
## 0  0.05  0.50  0.45  0.00  0.00  0.00  0.00  0.00
## 10 0.05  0.00  0.50  0.45  0.00  0.00  0.00  0.00
## 20 0.00  0.05  0.00  0.50  0.45  0.00  0.00  0.00
## 30 0.00  0.00  0.05  0.00  0.50  0.45  0.00  0.00
## 40 0.00  0.00  0.00  0.05  0.00  0.50  0.45  0.00
## 50 0.00  0.00  0.00  0.00  0.05  0.00  0.50  0.45
## 60 0.00  0.00  0.00  0.00  0.00  0.05  0.00  0.95
## 70 0.00  0.00  0.00  0.00  0.00  0.00  0.00  1.00
```

## Summary

1)  $\pi : SBA$

```
print(pi_speed)
```

```
##      normal  speed
## 0         0      1
## 10        0      1
## 20        0      1
## 30        0      1
## 40        0      1
## 50        0      1
## 60        0      1
## 70        0      1
```

```
print(pi_50)
```

```
##      normal  speed
## 0         0.5    0.5
## 10        0.5    0.5
## 20        0.5    0.5
## 30        0.5    0.5
## 40        0.5    0.5
## 50        0.5    0.5
## 60        0.5    0.5
## 70        0.5    0.5
```

2)  $R^\pi : SBR$

```
print(reward_fn(pi_speed))
```

```
## [[-1.5]
##  [-1.5]
##  [-1.5]
##  [-1.5]
##  [-0.5]
##  [-1.5]
##  [-1.5]
##  [ 0.  ]]
```

```
print(reward_fn(pi_50))
```

```
## [[-1.25]
##  [-1.25]
##  [-1.25]
##  [-1.25]
##  [-0.25]
##  [-1.25]
##  [-1.25]
##  [ 0.  ]]
```

3)  $P^\pi$  : SÖABS

```
print(transition(pi_speed, states=states, P_normal=P_normal, P_speed=P_speed))
```

```
##      0    10    20    30    40    50    60    70
## 0    0.1  0.0  0.9  0.0  0.0  0.0  0.0  0.0
## 10   0.1  0.0  0.0  0.9  0.0  0.0  0.0  0.0
## 20   0.0  0.1  0.0  0.0  0.9  0.0  0.0  0.0
## 30   0.0  0.0  0.1  0.0  0.0  0.9  0.0  0.0
## 40   0.0  0.0  0.0  0.1  0.0  0.0  0.9  0.0
## 50   0.0  0.0  0.0  0.0  0.1  0.0  0.0  0.9
## 60   0.0  0.0  0.0  0.0  0.0  0.1  0.0  0.9
## 70   0.0  0.0  0.0  0.0  0.0  0.0  0.0  1.0
```

```
print(transition(pi_50, states=states, P_normal=P_normal, P_speed=P_speed))
```

```
##      0    10    20    30    40    50    60    70
## 0    0.05 0.50 0.45 0.00 0.00 0.00 0.00 0.00
## 10   0.05 0.00 0.50 0.45 0.00 0.00 0.00 0.00
## 20   0.00 0.05 0.00 0.50 0.45 0.00 0.00 0.00
## 30   0.00 0.00 0.05 0.00 0.50 0.45 0.00 0.00
## 40   0.00 0.00 0.00 0.05 0.00 0.50 0.45 0.00
## 50   0.00 0.00 0.00 0.00 0.05 0.00 0.50 0.45
## 60   0.00 0.00 0.00 0.00 0.00 0.05 0.00 0.95
## 70   0.00 0.00 0.00 0.00 0.00 0.00 0.00 1.00
```

## Final Implementation

```
def policy_eval(given_pi):
    R=reward_fn(given_pi)
    P=transition(given_pi, states=states, P_normal=P_normal, P_speed=P_speed)

    gamma=1.0
    epsilon=10**(-8)

    v_old=np.repeat(0,8).reshape(8,1)
    v_new=R+np.dot(gamma*P, v_old)

    while np.max(np.abs(v_new-v_old))>epsilon:
        v_old=v_new
        v_new=R+np.dot(gamma*P,v_old)

    return v_new.T
print(policy_eval(pi_speed))
```

```
## [[-5.80592905 -5.2087811 -4.13926239 -3.47576467 -2.35376031 -1.73537603
##   -1.6735376  0.          ]]
```

```
print(policy_eval(pi_50))
```

```
## [[-5.96923786 -5.13359222 -4.11995525 -3.38922824 -2.04147003 -2.02776769
##   -1.35138838  0.          ]]
```