

Lecture E2. MDP with Model 2

Sim, Min Kyu, Ph.D., mksim@seoultech.ac.kr



서울과학기술대학교 데이터사이언스학과

- 1 I. Recap
- 2 II. Policy improvement
- 3 III. Policy iteration

I. Recap

policy_eval()

```

✓ gamma <- 1
✓ states <- as.character(seq(0, 70, 10))
✓ P_normal <- matrix(c(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),
                    nrow = 8, ncol = 8, byrow = TRUE,
                    dimnames = list(states, states))
✓ P_speed <- matrix(c(.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),
                   nrow = 8, ncol = 8, byrow = TRUE,
                   dimnames = list(states, states))

```

```

✓ transition <- function(given_pi,
                        states, P_normal, P_speed) {
  P_out <- array(0,
                dim = c(length(states), length(states)),
                dimnames = list(states, states))
  for (s in states) {
    action_dist <- given_pi[s,]
    P <- action_dist["normal"]*P_normal +
          action_dist["speed"]*P_speed
    P_out[s,] <- P[s,]
  }
  return(P_out)
}
✓ R_s_a <- matrix(
  c( -1,  -1,  -1,  -1,  0.0,  -1,  -1,  0,
     -1.5,-1.5,-1.5,-1.5,-0.5,-1.5,-1.5, 0),
  nrow = length(states), ncol = 2, byrow = FALSE,
  dimnames = list(states, c("normal", "speed")))
✓ reward_fn <- function(given_pi, R_s_a) {
  R_pi <- rowSums(given_pi*R_s_a)
  return(R_pi)
}

```

```

policy_eval <- function(given_pi) {
  R <- reward_fn(given_pi, R_s_a = R_s_a)
  P <- transition(given_pi, states = states, P_normal = P_normal, P_speed = P_speed)
  gamma <- 1.0
  epsilon <- 10^(-8)
  v_old <- array(rep(0,8), dim=c(8,1))
  v_new <- R + gamma*P%*%v_old
  while (max(abs(v_new-v_old)) > epsilon) {
    v_old <- v_new
    v_new <- R + gamma*P%*%v_old
  }
  return(v_new)
}

pi_speed <- cbind(rep(0,length(states)), rep(1,length(states)))
rownames(pi_speed) <- states; colnames(pi_speed) <- c("normal", "speed")
t(policy_eval(pi_speed))

```

```

##           0          10          20          30          40          50          60 70
## [1,] -5.805929 -5.208781 -4.139262 -3.475765 -2.35376 -1.735376 -1.673538 0

```

```

pi_50 <- cbind(rep(0.5,length(states)), rep(0.5,length(states)))
rownames(pi_50) <- states; colnames(pi_50) <- c("normal", "speed")
t(policy_eval(pi_50))

```

```

##           0          10          20          30          40          50          60 70
## [1,] -5.969238 -5.133592 -4.119955 -3.389228 -2.04147 -2.027768 -1.351388 0

```

Major components of approaching MDP

- ① **(policy evaluation)** We need to be able to evaluate $V^\pi(s)$ for a fixed π . This is called policy evaluation. This is also called as prediction in reinforcement learning.
 - ② **(optimal value function)** We want to be able to evaluate $V^{\pi^*}(s)$ where π^* is the *optimal policy*. The quantity, $V^{\pi^*}(s)$, is optimal policy's value function, or called shortly as *optimal value function*.
 - ③ **(optimal policy)** We want to find the *optimal policy* π^* . This is also called as *control* in reinforcement learning
- Check your reasoning why the followings are possible.
 - Optimal policy first: (optimal policy) + (policy evaluation) → (optimal value function)
 - Optimal value function first: (optimal value function) → (optimal policy)
 - This note will discuss
 - (policy evaluation) + series of (policy improvement) → (optimal policy)



II. Policy improvement

Development

- Remind that we have a Bellman's equation for MDP as follows.

$$\underbrace{V^\pi(s)}_{t \sim} = \underbrace{R^\pi(s)}_t + \gamma \underbrace{\sum_{\forall s'} P_{ss'}^\pi V^\pi(s')}_{t+1 \sim} \quad (\text{E1, p18})$$

- It means that, given a π , its value is determined by immediate reward plus the discounted sum of future reward.

- In this light, let's try to criticize the π_{speed} .

`t(policy_eval(pi_speed))`

##	0	10	20	30	40	50	60	70	
##	[1,]	-5.805929	-5.208781	-4.139262	-3.475765	-2.35376	-1.735376	-1.673538	0

Handwritten notes on the table:

- At state 40: $R + \gamma V$
 - i) normal: π
 - $-1 + \gamma(-2.35) = -3.35$
 - ii) speed: π
 - $-1.5 + \gamma(-1.73) + 0.1 \times (-4.14) = -3.48$
- At state 60: $V^{\pi_{\text{speed}}}(60) = -1.67$
- Equation: $= R^\pi(60) + \gamma P V$

- From the state 60, current policy gives the estimate for the state-value function of -1.6735376.
- We know that switching to *normal mode* at state 60 is better alternative than current action of *speed mode*. Because it guarantess the arrival to the state 70 with additional energy spending of 1.0.
- How would you express this fact in a mathematical form?

- Under the current policy's (π_{speed}) value function, on the state 60,

- Choosing normal mode gives

$$R + \gamma PV = \underline{-1.0} + \underline{1.0} \cdot \underline{0} = \underline{-1.0}$$

\downarrow $a = \text{normal}$ \downarrow $\pi = \pi_{\text{speed}}$

- Choosing speed mode gives

$$R + \gamma PV = \underline{-1.5} + (\underline{0.9} \cdot \underline{0} + \underline{0.1} \cdot \underline{-1.74}) = \underline{-1.674}$$

\downarrow $\pi = \pi_{\text{speed}}$ \downarrow $\pi = \pi_{\text{speed}}$

- This, π_{speed} should modify its action on the state 60.

- You just improved the current policy π_{speed} for the state 60!

- This should be checked for all states as well as the state 60.

- This completes **policy improvement**.

- Formally, policy improvement implies the following task of replacement:

$$\underline{\pi^{new}(s)} \leftarrow \underline{\operatorname{argmax}_{a \in \mathcal{A}}} \left[R(s, \underline{a}) + \gamma \sum_{s'} P_{ss'}^a V^{\underline{\pi^{old}}}(s') \right], \text{ for all } s$$

$\pi^{old}(30) = \text{sp.}$

$$\pi^{new}_{(30)} = \operatorname{argmax}_{a \in \mathcal{A}}$$

$$\begin{bmatrix} -1 + (2 \cdot 35) \\ -1.5 + (0.9 \times -1.74 + 0.1 \times -4.14) \end{bmatrix}$$

$$\text{for } s=30$$

$$\pi^{new}(s) \leftarrow \underset{a \in \mathcal{A}}{\operatorname{argmax}} \left[R(s, a) + \gamma \sum_{\forall s'} P_{ss'}^a V^{\pi^{old}}(s') \right], \text{ for all } s$$

- The term in the RHS, $R(s, \overset{t}{a}) + \gamma \sum_{\forall s'} \overset{t+1}{P_{ss'}^a} V^{\pi^{old}}(s')$, implies **[an expected return of starting from state s choosing an action a for this time step only, then following the policy π afterwards.]**
- How is this quantity different from $V^{\pi}(s)$?

$V(s)$: state-value fn
 $q(a, s)$: ~~state~~-action-value fn

- ✓ • The RHS makes an improvement using current policy π , by varying only the action in this time step.
- Formally, $q(s, a)$ is called action-value function, also famously known as Q-function.

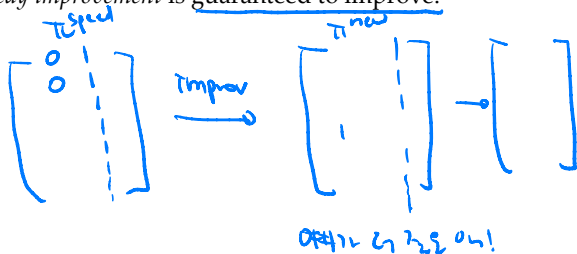
$$\begin{aligned} q^{\pi}(s, a) &:= \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a] \\ &= R(s, a) + \gamma \sum_{\forall s'} P_{ss'}^a V^{\pi^{old}}(s') \end{aligned}$$

$$\pi^{new}(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} \left[R(s, a) + \gamma \sum_{\forall s'} P_{ss'}^a V^{\pi^{old}}(s') \right], \text{ for all } s$$

- Using this new notation of $q(s, a)$, the policy improvement can be written as

$$\pi^{new}(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} q^{\pi^{old}}(s, a), \text{ for all } s$$

- The improvement is called *greedy improvement* since it involves a myopic digression from the current policy, in a way that an action only on this time step is revised.
- It can be proved that *greedy improvement* is guaranteed to improve.



Implementation

$$\pi^{new}(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} \left[R(s, a) + \gamma \sum_{s'} \mathbf{P}_{ss'}^a V^{\pi^{old}}(s') \right], \text{ for all } s$$

```
V_old <- policy_eval(pi_speed)
pi_old <- pi_speed
q_s_a <- R_s_a +
  cbind(gamma*P_normal%%V_old,
        gamma*P_speed%%V_old)
q_s_a
```

	normal	speed
## 0	-5.208781	5.805929
## 10	-5.139262	-5.208781
## 20	-4.475765	-4.139262
## 30	-3.353760	-3.475765
## 40	-1.735376	-2.353760
## 50	-2.673538	-1.735376
## 60	-1.000000	-1.673538
## 70	0.000000	0.000000

1 step improvement

```
pi_new_vec <- apply(q_s_a, 1, which.max)
pi_new <- array(0, dim = dim(pi_old),
               dimnames = dimnames(pi_old))
for (i in 1:length(pi_new_vec)) {
  pi_new[i, pi_new_vec[i]] <- 1
}
pi_new
```

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

[2 1 2 1 1 2 1]

- `policy_improve()`

```
policy_improve <- function(
  V_old,
  pi_old = pi_old, R_s_a = R_s_a, gamma = gamma,
  P_normal = P_normal, P_speed = P_speed) {

  q_s_a <- R_s_a + cbind(gamma*P_normal**V_old,
                        gamma*P_speed**V_old)

  pi_new_vec <- apply(q_s_a, 1, which.max)
  pi_new <- array(0, dim = dim(pi_old),
                 dimnames = dimnames(pi_old))

  for (i in 1:length(pi_new_vec)) {
    pi_new[i, pi_new_vec[i]] <- 1
  }
  return(pi_new)
}
```

code reuse

- One step improvement from π^{speed}

```
pi_old <- pi_speed ✓
V_old <- policy_eval(pi_old) ✓
pi_new <- policy_improve(V_old, ✓
  pi_old = pi_old, R_s_a = R_s_a, gamma = gamma,
  P_normal = P_normal, P_speed = P_speed)
```

pi_old

##	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

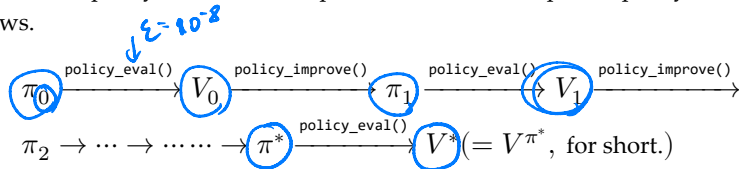
pi_new

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

III. Policy iteration

Discussion

- Given a policy π , policy_eval() evaluates its state-value function.
- Using the estimate of state-value function, policy_improve() improves the policy to the better one.
- If this process is iterated, then it is guaranteed to reach optimal policy.
- In other words, policy iteration is the process to reach the optimal policy described as follows.



- The iteration process terminates when π_i does not change any more, i.e. $\pi_i = \pi_{i+1}$.
- Note that policy evaluation is an approximate algorithm. For policy iteration purpose, policy evaluation cannot be, (and doesn't have to be as well), perfect.

Try do it over and over until no change - from π^{speed}

● Step 0

```
pi_old <- pi_speed  
pi_old
```

```
##      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
```

● Step 1

```
V_old <- policy_eval(pi_old)  
pi_new <- policy_improve(V_old,  
  pi_old = pi_old, R_s_a = R_s_a, gamma = gamma,  
  P_normal = P_normal, P_speed = P_speed)  
pi_old <- pi_new  
pi_old ✓
```

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



● Step 2

```
V_old <- policy_eval(pi_old)
pi_new <- policy_improve(V_old,
  pi_old = pi_old, R_s_a = R_s_a, gamma = gamma,
  P_normal = P_normal, P_speed = P_speed)
pi_old <- pi_new
pi_old
```

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



● Step 3

```
V_old <- policy_eval(pi_old)
pi_new <- policy_improve(V_old,
  pi_old = pi_old, R_s_a = R_s_a, gamma = gamma,
  P_normal = P_normal, P_speed = P_speed)
pi_old <- pi_new
pi_old
```

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

Policy iteration process (from π^{speed})

- Now we are ready to implement whole process as a single code block.

```
pi_old <- pi_speed
cnt <- 0
repeat{ # do-while in R
  print(paste0(cnt, "-th iteration"))
  print(t(pi_old))
  V_old <- policy_eval(pi_old)
  pi_new <- policy_improve(V_old,
    pi_old = pi_old, R_s_a = R_s_a, gamma = gamma,
    P_normal = P_normal, P_speed = P_speed)
  if (all.equal(pi_new, pi_old)==TRUE) break
  pi_old <- pi_new
  cnt <- cnt + 1
}
print(policy_eval(pi_new))
```

```
## [1] "0-th iteration"
##      0 10 20 30 40 50 60 70
## normal 0 0 0 0 0 0 0 0
## speed  1 1 1 1 1 1 1 1
## [1] "1-th iteration"
##      0 10 20 30 40 50 60 70
## normal 0 1 0 1 1 0 1 1
## speed  1 0 1 0 0 1 0 0
## [1] "2-th iteration"
##      0 10 20 30 40 50 60 70
## normal 0 0 0 1 1 0 1 1
## speed  1 1 1 0 0 1 0 0
##      [,1]
## 0 -5.107744
## 10 -4.410774
## 20 -3.441077
## 30 -2.666667
## 40 -1.666667
## 50 -1.666667
## 60 -1.000000
## 70 0.000000
```

Handwritten calculations and annotations:

- For state 30: $-1 - 1.67$
- For state 40: $-1.5 - (0.9 \times 0 + 0.1 \times -1.67)$
- Arrows indicate the calculation of the Bellman optimality backup for each state.
- A large bracket on the left groups the states 30 through 70.

Policy iteration process (from π^{50})

- The process should work for other initial choice of π , albeit possibly different convergence rate.

```
pi_old <- pi_50
cnt <- 0
repeat{ # do-while in R
  print(paste0(cnt, "-th iteration"))
  print(t(pi_old))
  V_old <- policy_eval(pi_old)
  pi_new <- policy_improve(V_old,
    pi_old = pi_old, R_s_a = R_s_a, gamma = gamma,
    P_normal = P_normal, P_speed = P_speed)
  if (all.equal(pi_new, pi_old)==TRUE) break
  pi_old <- pi_new
  cnt <- cnt + 1
}
print(policy_eval(pi_new))
```

```
## [1] "0-th iteration"
##           0 10 20 30 40 50 60 70
## normal 0.5 0.5 0.5 0.5 0.5 0.5 0.5 0.5
## speed  0.5 0.5 0.5 0.5 0.5 0.5 0.5 0.5
## [1] "1-th iteration"
##           0 10 20 30 40 50 60 70
## normal 0   1  0  1  1  0  1  1
## speed  1   0  1  0  0  1  0  0
## [1] "2-th iteration"
##           0 10 20 30 40 50 60 70
## normal 0   0  0  1  1  0  1  1
## speed  1   1  1  0  0  1  0  0
##           [,1]
## 0  -5.107744
## 10 -4.410774
## 20 -3.441077
## 30 -2.666667
## 40 -1.666667
## 50 -1.666667
## 60 -1.000000
## 70  0.000000
```

Summary

- From a policy π , $q^\pi(s, a)$ implies an expected return of starting from the state s , choosing an action a for this time step only, then following the policy π afterwards.
- Policy improvement occurs by

$$\pi^{new}(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} q^{\pi^{old}}(s, a) \quad \checkmark$$

, or, equivalently,

$$\pi^{new}(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} \left[R(s, a) + \gamma \sum_{\forall s'} P_{ss'}^a V^{\pi^{old}}(s') \right], \text{ for all } s \quad \checkmark$$

- Policy iteration is the iterative process from an arbitrary policy, policy evaluation and policy improvement take places until the policy converges. It is guaranteed to be converges to the optimal policy.

"Success isn't permanent, and failure isn't fatal. - Mike Ditka"