

→ MDP에 대한 모든 정보를 알 때, 최적의 풀리시를 찾아나가는 방법.

## Lecture 3: Planning by Dynamic Programming

↪ 방법론

David Silver

# Outline

## 1 Introduction

## 2 Policy Evaluation

policy 가 정해졌을 때, 이를 따라가면 어떤 결과를 얻는지  
value function 을 찾는 것 ( policy 평가 )

## 3 Policy Iteration

Iteration 한 방법으로 최적의 policy 를 찾는 것

## 4 Value Iteration

## 5 Extensions to Dynamic Programming

## 6 Contraction Mapping (X) Silver & 강희 안함

# What is Dynamic Programming?

**Dynamic** sequential or temporal component to the problem

**Programming** optimising a “program”, i.e. a policy

- c.f. linear programming

- A method for **solving complex problems**

복잡한 문제를 풀 방법은.

- By **breaking them down** into subproblems

큰 문제를 작게 쪼갬다.

- Solve the subproblems 쪼갬 것을 푼다.

- Combine solutions to subproblems 조합한다.

# Requirements for Dynamic Programming



Dynamic Programming is a very general solution method for problems which have two properties:

- **Optimal substructure**    최적의 해결책이 나와진 작은 문제에 적용이 되어야 한다
  - *Principle of optimality* applies
  - Optimal solution can be decomposed into subproblems
- **Overlapping subproblems**    많이 등장하니까
  - Subproblems recur many times    한 subproblem을 풀어서 저장해두면
  - Solutions can be cached and reused    다시 사용할 수 있다. (다시 안나기 때문)
- Markov decision processes satisfy both properties    MDP가 위 두 property를 만족한다.
  - Bellman equation gives recursive decomposition
  - Value function stores and reuses solutions

# Planning by Dynamic Programming

→ env 가 동작하는 것을 알고있음

- Dynamic programming assumes full knowledge of the MDP
- It is used for *planning* in an MDP
- For prediction:  $v^f$  을 찾음
  - Input: MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$  and policy  $\pi$
  - or: MRP  $\langle \mathcal{S}, \mathcal{P}^\pi, \mathcal{R}^\pi, \gamma \rangle$       ↳ 애가 종료지 내뽐어 optimal 인지 알면 안함
  - Output: value function  $v_\pi$
- Or for control: policy 를 찾음
  - Input: MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
  - Output: optimal value function  $v_*$
  - and: optimal policy  $\pi_*$

# Other Applications of Dynamic Programming

Dynamic programming is used to solve many other problems, e.g.

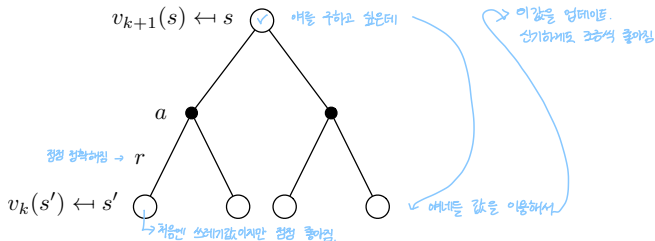
- Scheduling algorithms
- String algorithms (e.g. sequence alignment)
- Graph algorithms (e.g. shortest path algorithms)
- Graphical models (e.g. Viterbi algorithm)
- Bioinformatics (e.g. lattice models)

# Iterative Policy Evaluation 폴리시를 평가하자~

- 이 policy 를 따라갔을 때 return을 얼마받나 = policy Evaluation = value function
- prediction 문제이다

- Problem: evaluate a given policy  $\pi$
- Solution: iterative application of Bellman expectation backup
- $v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_\pi$  증적 목표 을 반복 실행 → 값만 저장하는 것
- Using *synchronous* backups, 한 단계에 한 개씩
  - At each iteration  $k + 1$
  - For all states  $s \in \mathcal{S}$
  - Update  $v_{k+1}(s)$  from  $v_k(s')$  모든 state 를 한 번씩 업데이트 뒤에 값  $v_k(s')$  을 이용해서  $v_{k+1}(s)$  을 업데이트
  - where  $s'$  is a successor state of  $s$
- We will discuss *asynchronous* backups later
- Convergence to  $v_\pi$  will be proven at the end of the lecture

# Iterative Policy Evaluation (2)



**Bellman expectation equation**

$$v_{k+1}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_k(s') \right)$$

$$\mathbf{v}^{k+1} = \mathcal{R}^\pi + \gamma \mathcal{P}^\pi \mathbf{v}^k$$

매 Iterative 마다 위 식을 이용하여 모든  $s$  를 업데이트



# Evaluating a Random Policy in the Small Gridworld

prediction 문제이기 때문에

policy 와 MDP 가 주어진 상태



(+w/e/n) actions

\* S가 없는데 S로 이동하면  
현재 있던 자리에 있음



Q. 1칸에서 14칸으로 갈 확률은? 1에서 3칸으로 갈 확률은?  
A. 확률다. 연속으로  $\frac{1}{4}$  이겠지만, 명령 들수도 있으니까.

$r = -1$   
on all transitions

다시 한번 문제 정리

MDP와 Policy 가 있는 상태. (prediction 문제)  
value를 찾는다.  
value는 한 상태에서 어떤 액션이 좋을지이다.  
(한번 실행하면 다음 -2가 된다)  
value는 state의 기대값이다. state에  
있는 액션이 주는 기대값과 그에 따른  
discounting이 있다.

↪ 이제 MDP 재

- Undiscounted episodic MDP ( $\gamma = 1$ )  $\hookrightarrow$  SA PR  $\gamma$  를 다 알고있으니까
- Nonterminal states 1, ..., 14
- One terminal state (shown twice as shaded squares)
- Actions leading out of the grid leave state unchanged
- Reward is  $-1$  until the terminal state is reached
- Agent follows uniform random policy 4가지 방향에 각 0.25 확률로 이동

$$\pi(n|\cdot) = \pi(e|\cdot) = \pi(s|\cdot) = \pi(w|\cdot) = 0.25$$

$\uparrow \uparrow$   
 $a \quad S$

# Iterative Policy Evaluation in Small Gridworld

일단 state 칸만큼 테이블을 만들어놓음

$V_k$  for the  
Random Policy

$k = 0$   
(0으로 초기화)

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

왜 이렇게 되냐? Bellman equation에 따라서!

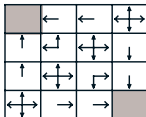
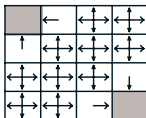
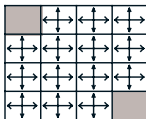
$k = 1$

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

$k = 2$

0.0	-1.7	-2.0	-2.0
-1.7	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.7
-2.0	-2.0	-1.7	0.0

Greedy Policy  
w.r.t.  $V_k$



random  
policy

$$V^{k+1} = \beta \pi + \gamma P^\pi V^k$$

maybe 1      배 곱셈

계산

$$V^1 = -1 + \gamma \left( \frac{1}{4} \cdot 0 + \frac{1}{4} \cdot 0 + \dots + \frac{1}{4} \cdot 0 \right)$$

$$V^2 = -1 + \gamma \left( \frac{1}{4} \cdot (-1) + \dots + \frac{1}{4} \cdot (-1) \right)$$

계산

$$V^2 = -1 + \gamma \left( \frac{1}{4} \cdot (-1) + \frac{1}{4} \cdot (0) + \frac{1}{4} \cdot (-1) + \frac{1}{4} \cdot (-1) \right)$$

무슨 값인가? 계산한거 가져와...

$$= -1 - \frac{3}{4} = -1.75$$

반올림한거임

계산한게 맞는지 확인하는 과정.  
값과 인덱스가 일치한것 같아...

사실 우리는 굉장히 다양한  
policy를 평가할 것인데,  
이 평가된 value 안에서  
greedy 하게 움직이면  
결국엔 optimal policy가  
찾아올 수 있다.

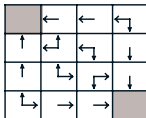
평가만 했을때!  
optimal  
policy  
가 생겼다!

# Iterative Policy Evaluation in Small Gridworld (2)

이 문제에서는  $k=3$ 까지만 해도 Optimal policy 가 나온다

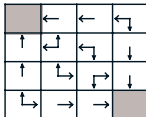
$k = 3$

0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0



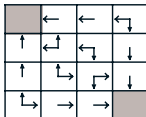
$k = 10$

0.0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0.0



$k = \infty$

0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0



optimal policy

~~~~~ Evaluate  $\rightarrow$  Improve  $\rightarrow$  Evaluate  $\rightarrow$  Improve ...

## How to Improve a Policy

- Given a policy  $\pi$

- Evaluate the policy  $\pi$     value function 을 찾는다.

$$v_{\pi}(s) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s]$$

- Improve the policy by acting greedily with respect to  $v_{\pi}$

찾은  $v_{\pi}$ 에 대해  $\pi'$ 을 만든다.

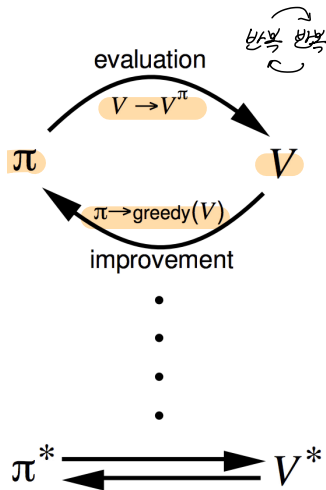
$$\pi' = \text{greedy}(v_{\pi})$$

방금 예제는 너무 간단해서  $\pi'$ 가 바로  $\pi^*$ 이 됐다. 보통은 여러번 해야  $\pi^*$ 에 도달

- In Small Gridworld improved policy was optimal,  $\pi' = \pi^*$
- In general, need more iterations of improvement / evaluation
- But this process of policy iteration always converges to  $\pi^*$

The diagram shows a series of points connected by arrows, representing an iterative improvement process. The points are arranged in a zig-zag pattern, moving from a 'starting' point towards a target point labeled  $V^*$  and  $\pi^*$ . The process is labeled 'improve' and  $V = V \pi$ . The target point is also labeled  $\pi = \text{greedy}(V)$ .

- Policy improvement Generate  $\pi' \geq \pi$
- Greedy policy improvement



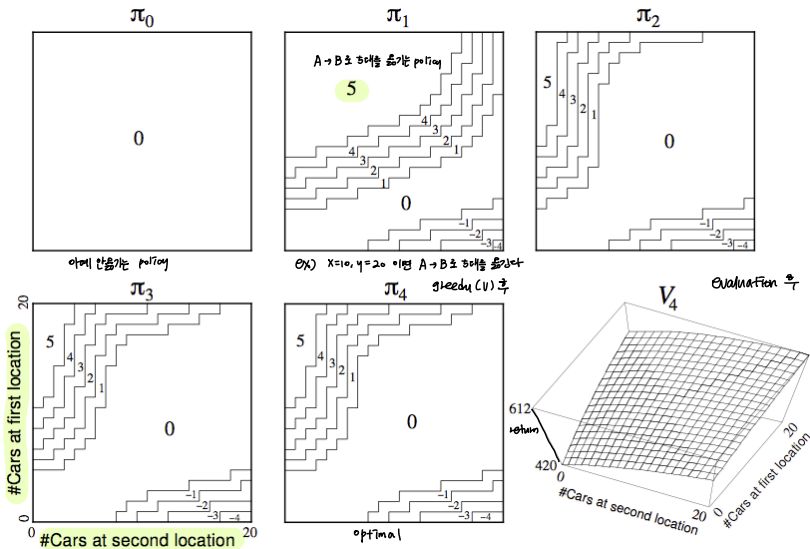
# Jack's Car Rental

렌트해주는 loc 2개.



- States: Two locations, **maximum of 20 cars** at each  
두 개의 렌트장이 있고, 한 곳에는 최대 20개의 차가 존재할 수 있다
- Actions: Move up to 5 cars between locations overnight  
A 와 B 간 차 이동이 가능하다.
- Reward: \$10 for each car rented (must be available)
- Transitions: Cars returned and requested randomly
  - Poisson distribution,  $n$  returns/requests with prob  $\frac{\lambda^n}{n!} e^{-\lambda}$   
정해진 일주일동안 차가 되돌아오는 확률
  - 1st location: average requests = 3, average returns = 3  
평균 하루에 3번 렌트 요청      평균 하루 3번 리턴
  - 2nd location: average requests = 4, average returns = 2

# Policy Iteration in Jack's Car Rental



# Policy Improvement

여기서 증명하고자 하는 것 : 정말 policy improvement 대로 하면 더 나은 policy 가 되는가?

결론 : 좋다!

- Consider a **deterministic policy**,  $a = \pi(s)$  어느 state 에 어떤 선택한 action 을 하는 policy
- We can **improve the policy by acting greedily** Q에 대해 greedy 하게 움직임

$$\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} q_{\pi}(s, a)$$

방법에서 볼 때,  
같은 S에서의 max 값보다  
적거나 같을 수 밖에

S로부터 pi를 따라가거나  
S에서 pi가 골라준 a를  
따라가거나 그 둘은 pi를  
따라가거나 같아

- This improves the value from any state  $s$  over one step,

(one step 이 때까진)

$$q_{\pi}(s, \pi'(s)) = \max_{a \in \mathcal{A}} q_{\pi}(s, a) \geq q_{\pi}(s, \pi(s)) = v_{\pi}(s)$$

greedy policy가 바가 될수 있는 것 중에  
가능한 것은 있을 것이기에

action value - function

S에 있을 때 바로 a를 고르고,  
 $q_{\pi}(s, a)$ 로  $v_{\pi}(s)$ 와 같다. (Q의 최대값 당면)

pi를 따를 때 따르는 V

- It therefore improves the value function,  $v_{\pi'}(s) \geq v_{\pi}(s)$

첫 step만 pi'를 따르고 다음 pi. by definition

$$v_{\pi}(s) \leq q_{\pi}(s, \pi'(s)) = \mathbb{E}_{\pi'} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s]$$

정확히의 원리나 관련 없다.

$$\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s]$$

$$\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 q_{\pi}(S_{t+2}, \pi'(S_{t+2})) \mid S_t = s]$$

$$\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \dots \mid S_t = s] = v_{\pi'}(s)$$

$v_{\pi} \rightarrow \pi$ 로  
변환시키면 공식  
을 생각해보니.

아니,,  
옳습니다.

(2번 반복적으로 적용해서  
필수 더 나은 결론.)



## Policy Improvement (2) 계속 나아지면 수렴 포인트는 optimal이다

- If improvements stop,

$$q_{\pi}(s, \pi'(s)) = \max_{a \in \mathcal{A}} q_{\pi}(s, a) = q_{\pi}(s, \pi(s)) = v_{\pi}(s)$$

- Then the Bellman optimality equation has been satisfied

논리 : 위 식이 만족되므로 optimal policy 다

$$v_{\pi}(s) = \max_{a \in \mathcal{A}} q_{\pi}(s, a)$$

- Therefore  $v_{\pi}(s) = v_{*}(s)$  for all  $s \in \mathcal{S}$
- so  $\pi$  is an optimal policy

# Modified Policy Iteration

- Does **policy evaluation** need to **converge to  $v_\pi$** ?  $v_\pi$ 로 수렴할 때까지 해야되나
- Or should we introduce a stopping condition **스탑을 만들어야 하나**
  - e.g.  $\epsilon$ -convergence of value function OR
- Or simply stop after  $k$  iterations of iterative policy evaluation?  **$k$ 번만?**
- For example, in the small gridworld  $k = 3$  was sufficient to achieve optimal policy  **$k$ 번만 해도 충분하다.**
- Why not update policy every iteration? i.e. stop after  $k = 1$ 
  - This is equivalent to *value iteration* (next section) **해도 된다.**

# Generalised Policy Iteration

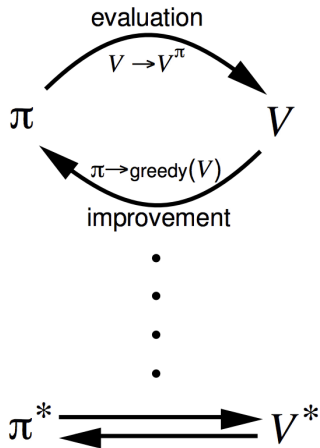


**Policy evaluation** Estimate  $v_\pi$

**Any** policy evaluation algorithm

**Policy improvement** Generate  $\pi' \geq \pi$

**Any** policy improvement algorithm



# Principle of Optimality

평요람이 이해하지 못해 PASS!

Any optimal policy can be subdivided into two components:

- An optimal first action  $A_*$
- Followed by an optimal policy from successor state  $S'$

## Theorem (Principle of Optimality)

*A policy  $\pi(a|s)$  achieves the optimal value from state  $s$ ,  $v_\pi(s) = v_*(s)$ , if and only if*

- *For any state  $s'$  reachable from  $s$*
- *$\pi$  achieves the optimal value from state  $s'$ ,  $v_\pi(s') = v_*(s')$*

# Deterministic Value Iteration

\* 앞과 다른 점: policy 가 없다.

- If we know the solution to subproblems  $v_*(s')$  이걸 알면
- Then solution  $v_*(s)$  can be found by one-step lookahead 이걸 구할수있다

< Bellman Optimality Equation >

$$v_*(s) \leftarrow \max_{a \in \mathcal{A}} \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s')$$

- The idea of value iteration is to apply these updates iteratively
- Intuition: start with final rewards and work backwards 마지막에서부터 계산.  
직관적
- Still works with loopy, stochastic MDPs  $\rightarrow$  결국적으로 optimal 에 수렴함.  
여기도 적용됨

## Example: Shortest Path

policy 가 없음!

|   |  |  |  |
|---|--|--|--|
| g |  |  |  |
|   |  |  |  |
|   |  |  |  |
|   |  |  |  |

Problem

|   |   |   |   |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

 $V_1$ 16개 값을 다 업데이트  
할건데 앞값의 식으로 계산.

|    |    |    |    |
|----|----|----|----|
| 0  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |

 $V_2$ 

|    |    |    |    |
|----|----|----|----|
| 0  | -1 | -2 | -2 |
| -1 | -2 | -2 | -2 |
| -2 | -2 | -2 | -2 |
| -2 | -2 | -2 | -2 |

 $V_3$ 

|    |    |    |    |
|----|----|----|----|
| 0  | -1 | -2 | -3 |
| -1 | -2 | -3 | -3 |
| -2 | -3 | -3 | -3 |
| -3 | -3 | -3 | -3 |

 $V_4$ 

|    |    |    |    |
|----|----|----|----|
| 0  | -1 | -2 | -3 |
| -1 | -2 | -3 | -4 |
| -2 | -3 | -4 | -4 |
| -3 | -4 | -4 | -4 |

 $V_5$ 

|    |    |    |    |
|----|----|----|----|
| 0  | -1 | -2 | -3 |
| -1 | -2 | -3 | -4 |
| -2 | -3 | -4 | -5 |
| -3 | -4 | -5 | -5 |

 $V_6$ 

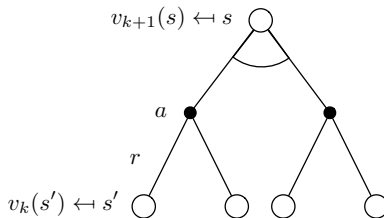
|    |    |    |    |
|----|----|----|----|
| 0  | -1 | -2 | -3 |
| -1 | -2 | -3 | -4 |
| -2 | -3 | -4 | -5 |
| -3 | -4 | -5 | -6 |

 $V_7$

# Value Iteration

- Problem: find optimal policy  $\pi$
- Solution: iterative application of Bellman optimality backup
- $v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_*$
- Using synchronous backups 각 반복마다 모든  $s$ 를 update
  - At each iteration  $k + 1$
  - For all states  $s \in \mathcal{S}$
  - Update  $v_{k+1}(s)$  from  $v_k(s')$
- Convergence to  $v_*$  will be proven later
- Unlike policy iteration, there is no explicit policy
- Intermediate value functions may not correspond to any policy  
중간값에 나오는 value function은 어떠한 policy도 아니다.

# Value Iteration (2)



$$v_{k+1}(s) = \max_{a \in \mathcal{A}} \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_k(s') \right)$$

$$\mathbf{v}_{k+1} = \max_{a \in \mathcal{A}} \mathbf{R}^a + \gamma \mathbf{P}^a \mathbf{v}_k \quad \leftarrow \text{Bellman}$$



## Example of Value Iteration in Practice

<http://www.cs.ubc.ca/~poole/demos/mdp/vi.html>

## Synchronous Dynamic Programming Algorithms

| Problem    | Bellman Equation                                         | Algorithm                   |
|------------|----------------------------------------------------------|-----------------------------|
| Prediction | Bellman Expectation Equation                             | Iterative Policy Evaluation |
| Control    | Bellman Expectation Equation + Greedy Policy Improvement | Policy Iteration            |
| Control    | Bellman Optimality Equation                              | Value Iteration             |

- Algorithms are based on state-value function  $v_{\pi}(s)$  or  $v_*(s)$
- Complexity  $O(mn^2)$  per iteration, for  $m$  actions and  $n$  states  
→ Complexity 가 크다
- Could also apply to action-value function  $q_{\pi}(s, a)$  or  $q_*(s, a)$
- Complexity  $O(m^2 n^2)$  per iteration

# Asynchronous Dynamic Programming

- DP methods described so far used *synchronous* backups
- i.e. all states are backed up *in parallel* 병렬
- *Asynchronous DP* backs up states individually, in any order 모든 상태가 아닌, 일부 / 개인만 backup.
- For each selected state, apply the appropriate backup
- Can significantly *reduce computation*
- Guaranteed to converge if *all states continue to be selected* 이게 보장되어야 함

# Asynchronous Dynamic Programming

Three simple ideas for asynchronous dynamic programming:

- *In-place* dynamic programming
- *Prioritised sweeping*
- *Real-time* dynamic programming

## In-Place Dynamic Programming

코딩 테크닉에 가깝네요.

- Synchronous value iteration stores two copies of value function

for all  $s$  in  $\mathcal{S}$ 테이블이 두 개 필요.  
old와 new

$$v_{new}(s) \leftarrow \max_{a \in \mathcal{A}} \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_{old}(s') \right)$$

$$v_{old} \leftarrow v_{new}$$

- In-place value iteration only stores one copy of value function

for all  $s$  in  $\mathcal{S}$ 테이블을 하나만,  
과거 정보 필요가 없이

$$v(s) \leftarrow \max_{a \in \mathcal{A}} \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v(s') \right)$$

# Prioritised Sweeping

우선순위를 두어서 업데이트. 벨만 에러를 편별.  
(가운뎃지)

- Use **magnitude of Bellman error** to guide state selection, e.g.

$$\left| \max_{a \in \mathcal{A}} \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v(s') \right) - v(s) \right|$$

- Backup the state with the largest remaining Bellman error
- Update Bellman error of affected states after each backup
- Requires knowledge of reverse dynamics (predecessor states)
- Can be implemented efficiently by maintaining a priority queue

# Real-Time Dynamic Programming

5가 많지만 agent가 보게 안함,  
agent가 볼때 각상태는 5만 보지 업데이트

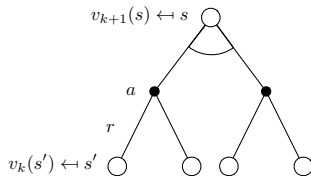
- Idea: only states that are relevant to agent
- Use agent's experience to guide the selection of states
- After each time-step  $S_t, A_t, R_{t+1}$
- Backup the state  $S_t$

$$v(S_t) \leftarrow \max_{a \in \mathcal{A}} \left( \mathcal{R}_{S_t}^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{S_t s'}^a v(s') \right)$$

# Full-Width Backups 여태까지 한 것.

- DP uses *full-width* backups
- For each backup (sync or async)
  - Every successor state and action is considered
  - Using knowledge of the MDP transitions and reward function
- DP is effective for *medium-sized* problems (millions of states)
- For large problems DP suffers *Bellman's curse of dimensionality*
  - Number of states  $n = |S|$  grows exponentially with number of state variables
- Even one backup can be too expensive

모든 갈 수 있는  $s$  가 고려됨.  
업데이트 시 참조하는  $s$  들.



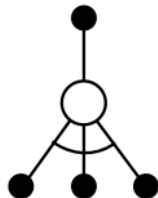
큰 문제에서 full-width는  $\times$  curse of dimensionality 문제



# Sample Backups

어디로 도착할지 몰라도, Backup 이 가능한게 장점!

- In subsequent lectures we will consider *sample backups*
- Using sample rewards and sample transitions  
 $\langle S, A, R, S' \rangle$
- Instead of reward function  $\mathcal{R}$  and transition dynamics  $\mathcal{P}$
- Advantages:  $\hookrightarrow$  우리 예제 Model-based 였대. free는 a한 어디까지모름
  - Model-free: no advance knowledge of MDP required
  - Breaks the curse of dimensionality through sampling
  - Cost of backup is constant, independent of  $n = |\mathcal{S}|$



# Approximate Dynamic Programming

silver

✓ 강의에서 다루지 않습니다!

- Approximate the value function
- Using a *function approximator*  $\hat{v}(s, \mathbf{w})$
- Apply dynamic programming to  $\hat{v}(\cdot, \mathbf{w})$
- e.g. Fitted Value Iteration repeats at each iteration  $k$ ,
  - Sample states  $\tilde{\mathcal{S}} \subseteq \mathcal{S}$
  - For each state  $s \in \tilde{\mathcal{S}}$ , estimate target value using Bellman optimality equation,

$$\tilde{v}_k(s) = \max_{a \in \mathcal{A}} \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \hat{v}(s', \mathbf{w}_k) \right)$$

- Train next value function  $\hat{v}(\cdot, \mathbf{w}_{k+1})$  using targets  $\{\langle s, \tilde{v}_k(s) \rangle\}$

## Some Technical Questions

- How do we know that value iteration converges to  $v_*$ ?
- Or that iterative policy evaluation converges to  $v_\pi$ ?
- And therefore that policy iteration converges to  $v_*$ ?
- Is the solution unique?
- How fast do these algorithms converge?
- These questions are resolved by *contraction mapping theorem*

# Value Function Space

- Consider the vector space  $\mathcal{V}$  over value functions
- There are  $|\mathcal{S}|$  dimensions
- Each point in this space fully specifies a value function  $v(s)$
- What does a Bellman backup do to points in this space?
- We will show that it brings value functions *closer*
- And therefore the backups must converge on a unique solution

## Value Function $\infty$ -Norm

- We will measure distance between state-value functions  $u$  and  $v$  by the  $\infty$ -norm
- i.e. the largest difference between state values,

$$\|u - v\|_{\infty} = \max_{s \in \mathcal{S}} |u(s) - v(s)|$$

# Bellman Expectation Backup is a Contraction

- Define the *Bellman expectation backup operator*  $T^\pi$ ,

$$T^\pi(v) = \mathcal{R}^\pi + \gamma \mathcal{P}^\pi v$$

- This operator is a  $\gamma$ -contraction, i.e. it makes value functions closer by at least  $\gamma$ ,

$$\begin{aligned} \|T^\pi(u) - T^\pi(v)\|_\infty &= \|(\mathcal{R}^\pi + \gamma \mathcal{P}^\pi u) - (\mathcal{R}^\pi + \gamma \mathcal{P}^\pi v)\|_\infty \\ &= \|\gamma \mathcal{P}^\pi(u - v)\|_\infty \\ &\leq \|\gamma \mathcal{P}^\pi\| \|u - v\|_\infty \\ &\leq \gamma \|u - v\|_\infty \end{aligned}$$

# Contraction Mapping Theorem

## Theorem (Contraction Mapping Theorem)

*For any metric space  $\mathcal{V}$  that is complete (i.e. closed) under an operator  $T(v)$ , where  $T$  is a  $\gamma$ -contraction,*

- *$T$  converges to a unique fixed point*
- *At a linear convergence rate of  $\gamma$*

# Convergence of Iter. Policy Evaluation and Policy Iteration

- The Bellman expectation operator  $T^\pi$  has a unique fixed point
- $v_\pi$  is a fixed point of  $T^\pi$  (by Bellman expectation equation)
- By contraction mapping theorem
- Iterative policy evaluation converges on  $v_\pi$
- Policy iteration converges on  $v_*$



# Bellman Optimality Backup is a Contraction

- Define the *Bellman optimality backup operator*  $T^*$ ,

$$T^*(v) = \max_{a \in \mathcal{A}} \mathcal{R}^a + \gamma \mathcal{P}^a v$$

- This operator is a  $\gamma$ -contraction, i.e. it makes value functions closer by at least  $\gamma$  (similar to previous proof)

$$\|T^*(u) - T^*(v)\|_\infty \leq \gamma \|u - v\|_\infty$$

# Convergence of Value Iteration

- The Bellman optimality operator  $T^*$  has a unique fixed point
- $v_*$  is a fixed point of  $T^*$  (by Bellman optimality equation)
- By contraction mapping theorem
- Value iteration converges on  $v_*$