



# TDDC17: Intro to Automated Planning

#### **Planning under Uncertainty**

Jonas Kvarnström

Artificial Intelligence and Integrated Computer Systems Division Department of Computer and Information Science Linköping University

jonas.kvarnstrom@liu.se - 2021

#### **Multiple Outcomes**



- Classical planning assumes we know outcomes in advance
  - State + action → unique resulting state
- Sometimes we <u>must</u> deal with <u>multiple outcomes</u>
  - Due to problems in execution
    - pick-up(object)
       Intended outcome: carrying(object) is true
       Unintended outcome: carrying(object) is false
  - Due to <u>random</u> but clearly <u>desirable / undesirable</u> outcomes
    - Toss a coin do I win?
  - Due to <u>random</u> outcomes with <u>unknown</u> long term effects
    - Do I end up in group A or B?
       No idea which one will turn out to be better for me

# Information, while planning



#### **Non-Deterministic Planning**

Model says: we end up in **<u>one of</u>** these states



- First "info dimension":
  - What <u>do</u> we know about action outcomes <u>when we create the plan</u>?

#### **Probabilistic Planning**

Model says: we end up in one of these states



#### Focus of this lecture!



### Probabilistic Planning: Defining the World as a Stochastic System

### **State Transition System**

#### • Classical planning: A <u>state transition system</u> $\Sigma = (S, A, \gamma)$

- $S = \{ s_0, s_1, \dots \}$ :
- $A = \{ a_0, a_1, \dots \}$ :
- $\gamma: S \times A \rightarrow 2^S$ :

- Finite set of **world states**
- Finite set of actions
- State transition function,

onkv@ida

specifying all "edges"



#### **Stochastic System**

#### **Probabilistic planning** uses a **stochastic system** $\Sigma = (S, A, P)$

- $S = \{ s_0, s_1, \dots \}$ :
- $A = \{ a_0, a_1, \dots \}$ :
- *P*(*s*, *a*, *s*'):

Finite set of **world states** 

Finite set of actions

Given that we are in s and execute *a*, the **probability** of ending up in s'

Replaces y

#### Planning Model says: we end up in one of these states Start here... A1 (0,1) (

# Stochastic System Example





- Some transitions are <u>deterministic</u>, some are <u>stochastic</u>
  - Trying to move from 12 to 13: You may end up at 15 instead (20% risk)
  - Trying to move from 11 to 14: You may stay where you are instead (50% risk)

# **Stochastic System: Planning and Execution**

#### • **Example**: In state **s1**, there are three possible actions

- The planner chooses the <u>action</u>
  - For example, move(l1,l4)
- Much later, when we actually execute this action:
  - We will find out the outcome (end up in s4 or s1), due to full observability
- But when we're planning:
  - Must try to prepare for both potential outcomes!



q

#### Fully Observable Probabilistic Planning: Policies and Histories

Important concepts, before we define the planning problem itself!

#### Policies; Example 1



- One type of formal plan structure: **Policy**  $\pi : S \to A$ 
  - Defining, <u>for each state</u>, which action to execute <u>whenever</u> we are there



Reaches s4 or s5, waits there infinitely many times

### Policy Example 2



• Example 2



Always reaches state s4, waits there infinitely many times

# Policy Example 3





Reaches state s4 with 100% probability "in the limit" (the more steps, the greater the probability)

#### **Policies and Histories**



- The <u>outcome</u> of sequentially executing a policy:
  - A <u>state sequence</u>  $h = \langle s_0, s_1, s_2, s_3, s_4, \dots \rangle$ , called a <u>history</u>
  - Infinite, since policies have no termination criterion
- For each policy, there can be many potential histories
  - Which one is the actual result? Gradually discovered at execution time!





Even if we know we start in s1: Two possible histories

- $h_1 = \langle s1, s2, s3, s4, s4, ... \rangle$  Reached s4, waits indefinitely
- $h_2 = \langle s1, s2, s5, s5 \dots \rangle$  Reached s5, waits indefinitely

How probable are these histories?

#### **Probabilities: Initial States, Transitions**

- Each policy has a **probability distribution** over **histories/outcomes** 
  - With known fixed initial state s<sub>0</sub>:

$$P(\langle s_0, s_1, s_2, s_3, \dots \rangle \mid \pi) = \prod_{i \ge 0} P(s_i, \pi(s_i), s_{i+1})$$
  
Probabilities for each required state transition

With unknown initial state:





• Two possible histories, if P(s1) = 1:

• 
$$h_1 = \langle s1, s2, s3, s4, s4, ... \rangle$$
  $-P(h_1 \mid \pi_1) = 1 \times 1 \times 0.8 \times 1 \times ... = 0.8$   
 $h_2 = \langle s1, s2, s5, s5 ... \rangle$   $-P(h_2 \mid \pi_1) = 1 \times 1 \times 0.2 \times 1 \times ... = 0.2$   
 $-P(h \mid \pi_1) = 1 \times 0 = 0$  for all other  $h$ 





•  $h_1 = \langle s1, s2, s3, s4, s4, ... \rangle$ 

 $P(h_1 \mid \pi_2) = 1 \times 1 \times 0.8 \times 1 \times ... = 0.8$  $h_3 = \langle s1, s2, s5, s4, s4, ... \rangle$   $P(h_3 \mid \pi_2) = 1 \times 1 \times 0.2 \times 1 \times ... = 0.2$  $P(h \mid \pi_2) = 1 \times 0$  for all other *h* 





### Objectives: Rewards for Indefinite Execution

# **Generalizating our Objectives**

- Policies allow indefinite execution
  - No predetermined termination criterion go on "forever"
    - def execute\_policy(π):
       while True:

```
s = sense_current_state()
```

```
a = \pi(s)
execute_action(a)
```

But without <u>termination</u>, there can't be <u>goal states</u>... So what is the objective? What is a <u>good</u> policy?

- Combination of:
  - **Cost function**  $c: (S, A) \to \mathbb{R}$ 
    - $c(s, a) = \text{cost of } \underline{\text{being in}}$  state s and  $\underline{\text{executing}}$  action a
  - **Reward function**  $R: (S, A, S) \to \mathbb{R}$ 
    - R(s, a, s') = Reward for being in s, executing a and actually ending up in s'

# **Example: Grid World**

#### Example: Grid World

- Actions: North, South, West, East, NorthWest, ..., TakeGold
  - Cost c(s, a) = 15 for all s, a
  - 90% chance: Go where you want
  - 10% risk: End up somewhere else
- <u>Rewards</u> for some transitions
  - R(s, a, s') = +100 for transitions when you take the gold in the top right cell
  - s = [top right, there is gold]
     a = TakeGold
    - s' = [top right, there is no gold]
- Danger in some cells
  - Try to go to the top right cell
  - R(s, a, s') = 0 usually
  - R(s, a, s') = -200
     <u>if</u> you accidentally end up in the danger cell



Important: States != locations Can't take the gold twice, can't gain infinite rewards



#### **Example: Tetris**



- In each "step", <u>a piece falls one row</u> and you execute <u>one action</u>
  - Guide the pieces left/right, rotate them, drop them
- If an action/step results in <u>filling a row</u>:
  - The line disappears
  - *R*(*old state, action, state with row removed*) = 100
- If an action/step results in <u>filling two rows</u>:
  - R(old state, action, state with 2 rows removed) = 400
- When a piece has <u>fallen all the way</u>:
  - A new *random* piece falls from the top
  - Model piece probabilities using P(s, a, s') according to (most types of) Tetris rules





#### **Example: Robot Navigation**

24 Pilowyuoj

- Example costs in robot navigation:
  - c(s, a) = 1 for each "horizontal" action
  - c(s, a) = 100 for each "vertical" action: Far away, difficult, ...
  - c(s, wait) = 1



# Example: Robot Navigation (2)

25 Piove

- Example rewards in robot navigation
  - Every time you **end up in s5**:
    - Negative reward maybe the robot is in our way
  - Every time you **end up in s4**:
    - Positive reward maybe it helps us



### Simplification



- To simplify formulas, include the cost in the reward!
  - Decrease each  $R(s_i, \pi(s_i), s_{i+1})$  by  $C(s_i, \pi(s_i))$



#### **Utility Functions**

How useful is an outcome to us?

# Total Rewards – In Advance?

- Given a policy  $\pi$ , what <u>will</u> our total rewards be?
  - Can't know in advance
    - Will I reach the goal or end up in the danger zone?
    - Which pieces will I get?







### Total Rewards – After Executing?

- Given a policy  $\pi$ ...
  - ...<u>and</u> an <u>outcome</u>, a <u>history</u> (*infinite* state sequence)  $h = \langle s_0, s_1, s_2, ... \rangle$  resulting from actually having executed  $\pi$ ...
- ...What <u>were</u> our total rewards?
  - Undiscounted utility of a history:

$$V(h|\pi) = \sum_{i \ge 0} R(s_i, \pi(s_i), s_{i+1})$$

- I was in  $s_0$ , executed  $\pi(s_0)$ , and ended up in  $s_1$  -- reward!
- I was in  $s_1$ , executed  $\pi(s_1)$ , and ended up in  $s_2$  -- reward!
- I was in  $s_2$ , executed  $\pi(s_2)$ , and ended up in  $s_3$  -- reward!
- ••••



#### **Discount Factors and Discounted Utility**

#### **Utility in a Context**



#### Policy = solution for <u>infinite</u> horizon

#### (Infinite actual execution)

#### Never ends – unrealistic

#### Indefinite execution

We will stop at some point (the universe will end), but we can't predict when

To find the best policy for *long term execution*: Consider the infinite case

#### **Infinite Undiscounted Utility**

- 32 Jonkwold
- If we use <u>undiscounted utility</u> for an <u>infinite history</u>:
  - $\pi_1$  **could** result in  $h_1$  = (s1, s2, s3, s4, s4, ... )
    - Stays at s4 forever, executing "wait" → <u>infinite</u> amount of rewards!
    - $V(h_1 \mid \pi_1) = (-100) + (-1) + (-100) + 100 + 100 + 100 + 100 + 100 + \dots$



# Infinite Undiscounted Utility (2)

#### What's the problem, if we "like" being in state s4?

- Can't distinguish between **different ways** of getting there!
  - $s1 \rightarrow s2 \rightarrow s3 \rightarrow s4$ :  $-201 + \infty = \infty$
  - $s1 \rightarrow s2 \rightarrow s1 \rightarrow s2 \rightarrow s3 \rightarrow s4$ :  $-401 + \infty = \infty$
  - Both appear equally good...
- Can't distinguish between infinite times **100** and and infinite times **1000**
- Even without infinity, we can't see the difference between rewards <u>now</u> and rewards in the <u>far future</u>





#### **Discounted Utility**

34 Star

- Solution: **Discounted utility** for a history
  - Introduce a **discount factor**,  $\gamma$ , with  $0 \le \gamma \le 1$
  - Let

$$V(h|\pi) = \sum_{i \ge 0} \frac{\gamma^i}{\gamma^i} R(s_i, \pi(s_i), s_{i+1})$$

- Distant rewards/costs have <u>less influence</u>
  - For example: 0.9, 0.81, 0.729, ...
- Discounted utility is <u>finite</u> as long as  $0 \le \gamma < 1$

Examples will use  $\gamma = 0.9$ 

*Only* to simplify formulas! Should choose carefully...



#### **Expected Utility of a Policy**

# **Expected Utility of a Policy**



 $V(h|\pi) = \sum_{i>0} \gamma^i R(s_i, \pi(s_i), s_{i+1})$ 

- We want to choose a **good policy** 
  - We know, for **<u>each</u>** history (outcome)  $h = \langle s_0, s_1, s_2, s_3, s_4, ... \rangle$  of a policy  $\pi$ :
    - The **probability** that the history will occur:  $P(h|\pi)$
    - The resulting actual discounted utility:
  - Using this, calculate the statistically <u>expected utility</u> (~"average" utility) for the entire <u>policy</u>:

$$E(\pi) = \sum_{h \in \{\text{all possible histories for }\pi\}} P(h|\pi)V(h|\pi)$$

• Or, the **expected utility** given that we **start execution in state s**:

$$E(\pi, s) = \sum_{h \in \{\text{all possible histories for }\pi\}} P(h \mid \pi, s_0 = s) V(h \mid \pi)$$

Strictly speaking, this is the expected <u>discounted</u> utility (but expected <u>undiscounted</u> utility is rarely used)
### **Remembering the Notation...**

Jonk@ida

- How to remember the notation?
  - V The actual Value. Value depends on rewards, and rewards depend on the state we actually ended up in, so we can only know the actual value if we know the actual outcome, which is represented as a history.
  - E the *Expected* value, in the statistical sense. Since it is only *expected*, it isn't certain yet, and that's because it takes into account *all possible outcomes* (histories) and the resulting rewards and probabilities.

### Example 1





### Example 2





Expected Utility of a Policy: Another Definition

# **Expected Utility: Example**

Consider a policy... 

Constantinenter

Conservation and a service of the se

In state A, we should execute the "green action", which might lead to:

Α

Conservation

Constant and a second second

C

Н

- B  $\rightarrow$  execute "blue"  $\rightarrow$  E, F or G
- C  $\rightarrow$  execute "red"  $\rightarrow$  H

B

F

G

Ε

• D  $\rightarrow$  execute green  $\rightarrow$  I, J or K



D

Κ

## **Expected Utility: History-based**



We calculated expected utilities based on histories



# **Expected Utility: Step by Step**

- Another computation method:
  - We want  $E(\pi, A)$ , and the selected action is  $\pi(A) = green$
  - What's the **probability** of outcome B?
  - What's the **reward** for this outcome? R(A, green, B)
  - How much more will I get after arriving in B?  $E(\pi, B)$ , by definition!
  - $\gamma E(\pi, B)$ How much is that **worth** to me **now**? Α What's the **probability** of outcome C? D B \*\*\*\*\* С Κ F Ε G 6.....

P(A, green, B)

Н

# Expected Utility: Step by Step (2)



- If π is a policy, then
  - $E(π,s) = \sum_{s' \in S} P(s, π(s), s') * (R(s, π(s), s') + γ E(π,s'))$
  - The expected utility of continuing to execute  $\pi$  after having reached s
  - Is the sum, for all possible states  $s' \in S$  that you might end up in (outcomes),

of the probability  $P(s, \pi(s), s')$  of actually ending up in that state given the action  $\pi(s)$  chosen by the policy, times

the reward you get for this transition

plus the discount factor times the expected utility  $E(\pi, s')$  of **continuing**  $\pi$  from the new state s'



### Example





### Trivial?



### Seems like a <u>trivial calculation</u>!

- Each expected utility depends on a few others...
   so you'd just keep computing until you get to the end?
- $E(\pi_2, s2) =$   $0.8(-1 + \gamma E(\pi_2, s3)) +$  $0.2(-1 + \gamma E(\pi_2, s5))$



### **Equation System**



If π is a policy, then
 E(π,s) = Σ<sub>s' ∈ S</sub> P(s, π(s), s') \* (R(s, π(s), s') + γ E(π,s'))

### This is an **equation system**: |S| equations, |S| variables!

Use standard solution methods...

# (Equation system example)

- Clarification: What do we mean by equation system?
  - Suppose you have 2 equations
    - x = 5y
    - y = x 20
  - Can't solve recursively:
    - Start with x = 5y; you want the value of y
    - We know y = x 20; we want the value of x, which is 5y, which is ...
  - Can solve as an equation system:
    - x = 5(x-20)
    - x = 5x 100
    - 0 = 4x 100
    - 100 = 4x
    - 25 = *x*

### Fully Observable Probabilistic Planning: Markov Decision Processes

### Markov Decision Processes

- Underlying world model:
- Plan representation:
- Goal representation:
- Solution:

#### Stochastic system

**Policy** – which action to perform in **<u>any</u>** state

#### **Reward function**

#### An **optimal policy**

#### Definition:

An optimal policy  $\pi^*$  <u>maximizes expected utility for all states</u>: For all states s and alternative policies  $\pi$ ,  $E(\pi^*, s) \ge E(\pi, s)$ 

#### Given a policy, we can compute expected utility. How do we <u>find</u> a policy <u>maximizing</u> it?



### Simplification



In many formulations of MDPs (and our robotic example), rewards <u>do not depend on the outcome s'</u>!

### Let's simplify the upcoming examples a bit...

$$E(\pi, s) = \sum_{s' \in S} P(s, \pi(s), s') \cdot (R(s, \pi(s), s') + \gamma E(\pi, s'))$$

$$\Rightarrow E(\pi, s) = \sum_{s' \in S} P(s, \pi(s), s') \cdot (R(s, \pi(s)) + \gamma E(\pi, s'))$$

$$\Rightarrow E(\pi, s) = R(s, \pi(s)) + \sum_{s' \in S} P(s, \pi(s), s') \cdot \gamma E(\pi, s')$$

Solving an MDP, step 1: Local modifications

### **Properties of Local Changes**

- Given an MDP, suppose that:
  - You already have some arbitrary (even random) policy  $\pi$
  - You select an arbitrary state  $s_k$ , make a **local change** to  $\pi(s_k)$ 
    - Example:  $\pi(s_k) = move(l1, l3) \rightarrow \pi(s_k) = move(l1, l4)$
  - The change turns out to be a local improvement
    - It **increases** the expected utility  $E(\pi, s_k)$  for **this** particular state  $s_k$
- Then the change <u>cannot</u> <u>decrease</u>  $E(\pi, s')$  for <u>any</u> s'!

### A local improvement for one state is always a global improvement



# Properties of Local Changes (2)



increase  $E(n, s_m)$ 



$$E(\pi, s_m) = R(s, \pi(s_m)) + \sum_{s' \in s} P(s, \pi(s_m), s') \cdot \gamma E(\pi, s')$$
  
All of these remain unchanged!  
$$E(\pi, s_k) \text{ may occur here } (s' = s_k),$$
  
but only **positively**:  
lncrease  $E(\pi, s_k) \Rightarrow$ 

# **Properties of Local Changes (3)**



- How many times can you make local improvements?
  - There is a **finite number** of possible **policies**,  $|S| \cdot |A|$ 
    - The number of states |S| is finite
    - The number of actions |A| is finite
  - You <u>can't</u> make the <u>same improvement</u> twice
    - Every improvement leads to strictly greater expected utility
    - Can't "loop around"
  - → After a finite number of local improvements, no more local improvement will be possible

## Properties of Local Changes (4)



### Also:

 Every global improvement <u>can be reached</u> through such local improvements (no need to first make the policy worse, then better)

### → We can <u>find optimal solutions</u> through <u>local</u> improvements

No need to "think globally"

### But how do we find a local improvement?

Remember, finding expected utilities required solving an expensive equation system...

### Is a Local Change an Improvement?

- To find out <u>if</u> a change is an improvement:
  - Take the current policy  $\pi$ , with an expected utility:

$$E(\pi, s) = R(s, \pi(s)) + \sum_{s' \in S} P(s, \pi(s), s') \cdot \gamma E(\pi, s')$$

- Define the Q function:
  - What is the expected utility

     if you <u>first execute action</u> a,
     regardless of what the policy says,
     but then continue executing the <u>old policy</u>
     for all other steps?
  - $Q(\pi, s, a) = R(s, a) + \sum_{s' \in S} P(s, a, s') \cdot \gamma E(\pi, s')$

If  $Q(\pi, s, a) > E(\pi, s)$ , then setting  $\pi(s) = a$  would be an *improvement* to  $\pi$ . We know this much without solving a full equation system... Just not how large the improvement is!

### Preliminaries 2: Example

- Example:  $E(\pi, s1)$ 
  - The expected utility of following  $\pi$
  - Starting in s1, beginning with move(l1,l2)
- $Q(\pi, s1, move(l1, l4))$ 
  - The expected utility of being in s1, first executing move(l1, l4), then following policy  $\pi$
  - Only used to quickly find improvements





### One Solution Method for MDPs: Policy Iteration

### **Policy Iteration**



- General idea:
  - Start out with an **initial policy**, maybe randomly chosen
  - Calculate and store the <u>expected utility</u> of executing that policy for each state
  - <u>Update</u> the policy by making a <u>local</u> decision <u>for each state</u>:
     "Which action should my <u>improved</u> policy choose in this state?"
    - Use the actions that *appear to be best* according to the Q function, based on the actual expected utility for the <u>current</u> policy
    - For every state s:

 $\pi'(s) := \arg\max_{a \in A} Q(\pi, s, a)$ 

Iterate until the policy no longer changes

But what if there was an <u>even better</u> choice, which we don't see now because of our single step modification (Q)?

That's OK: We still have an *improvement*, which cannot prevent *future* improvements in the next iteration

### **First Iteration**

# Policy Iteration 1: Initial Policy $\pi_1$





### Policy Iteration 2: Expected Utility for $\pi_1$

- Calculate expected utilities for the current policy  $\pi_1$ 
  - Simple: Chosen transitions are deterministic **and** return to the same state!
    - $E(\pi,s) = \frac{R(s,\pi(s))}{P(s,\pi(s))} + \gamma \sum_{s' \in S} P(s,\pi(s),s') E(\pi,s')$

• 
$$E(\pi 1, s1) =$$
 $R(s1, wait)$  $+ \gamma$  $E(\pi 1, s1)$  $= -1$  $+ 0.9$  $E(\pi 1, s1)$ •  $E(\pi 1, s2) =$  $R(s2, wait)$  $+ \gamma$  $E(\pi 1, s2)$  $= -1$  $+ 0.9$  $E(\pi 1, s2)$ •  $E(\pi 1, s3) =$  $R(s3, wait)$  $+ \gamma$  $E(\pi 1, s3)$  $= -1$  $+ 0.9$  $E(\pi 1, s3)$ •  $E(\pi 1, s4) =$  $R(s4, wait)$  $+ \gamma$  $E(\pi 1, s4)$  $= +100$  $+ 0.9$  $E(\pi 1, s4)$ •  $E(\pi 1, s5) =$  $R(s5, wait)$  $+ \gamma$  $E(\pi 1, s5)$  $= -100$  $+ 0.9$  $E(\pi 1, s5)$ 

- Simple equations to solve:
  - $0.1E(\pi 1, s1) = -1$
  - $0.1E(\pi 1,s2) = -1$
  - $0.1E(\pi 1, s3) = -1$
  - 0.1E(π1,s4) = +100
  - $0.1E(\pi 1, s5) = -100$

→ 
$$E(\pi 1, s2) = -10$$

→ 
$$E(π1,s3) = -10$$

→ 
$$E(π1,s4) = +1000$$

→ 
$$E(π1,s5) = -1000$$

#### **Given this policy** $\pi_1$ :

63

High rewards if we start in s4, high costs if we start in s5

# **Policy Iteration 3: Update 1a**





 But the values will yield good guidance to find <u>policy improvements</u>

# Policy Iteration 4: Update 1b



What is the best Iocal modification according to the <u>expected utilities</u> of the <u>current</u> policy?	$E(\pi_{1}, s1) = -10$ $E(\pi_{1}, s2) = -10$ $E(\pi_{1}, s3) = -10$ $E(\pi_{1}, s4) = +1000$ $E(\pi_{1}, s5) = -1000$			100
<ul> <li>For every state s</li> <li>Let π<sub>2</sub>(s) = argma</li> </ul>	: х <sub>а ∈ А</sub> Q(π1,s,a)			
<ul> <li>That is, find the a</li> <li>s2: wait move(l2,l1) move(l2,l3)</li> </ul>	ction <i>a</i> that max -1 + 0.9 * - -100 + 0.9 * - -1 + 0.9 * (0)	imize <mark>s R(s, a)</mark> + γ∑ 10 -10 0.8*–10+ 0.2*–1000)	$E_{s' \in S} P(s, a, s') E(\pi 1, s')$ = -10 = -109 = -188,2	

# Policy Iteration 5: Update 1c



What is the best <b>local</b> modificationaccording to the <b>expected utilities</b> of the <b>current</b> policy?	$(\pi_{1}, s1) = -10$ $(\pi_{1}, s2) = -10$ $(\pi_{1}, s3) = -10$ $(\pi_{1}, s4) = +1000$ $(\pi_{1}, s5) = -1000$			
<ul> <li>For every state s:</li> <li>Let π<sub>2</sub>(s) = argmax<sub>a</sub></li> </ul>	$_{\in A} Q(\pi 1, s, a)$			
• That is, find the action a that maximize $\frac{S}{S}R(s, a) + \gamma \sum_{s' \in S} P(s, a, s') E(\pi 1, s')$				
<ul><li>s3: wait</li></ul>	<b>-1</b> + 0.9 * ·	-10	= -10	
move(l3,l2) move(l3,l4)	-1 + 0.9 * - -100 + 0.9 * -	-10 +1000	= -10 = +800	
• s4: wait move(14.11)	+100 + 0.9 * +99 + 0.9 * -	+1000 -10	= +1000 = +90	
<ul><li>s5: wait</li></ul>	<mark>-100</mark> + 0.9 * -	-1000	= -1000	
move(15,12)	<mark>-101</mark> + 0.9 * -	-10	= -110	
move(15,14)	<mark>-200</mark> +0.9*-	+1000	= +700	

### **Policy Iteration 6: Second Policy**



Q-values based

on one modified

(can't decrease!)

action, then

following  $\pi_1$ 

#### • This results in a **new policy**

 $E(\pi_{1},s_{1}) = -10$ 

 $E(\pi 1, s2) = -10$ 

 $E(\pi 1, s3) = -10$ 

 $E(\pi 1, s4) = +1000$ 

 $E(\pi 1, s5) = -1000$ 

Now we have made use of earlier indications that s4 seems to be a good state

 $\pi_1 = \{(s1, wait), \}$ 

(s2, wait),

(s3, wait),

(s4, wait),

(s5, wait)

→ Try to go there from s1 / s3 / s5!

No change in s2 yet...



>=+444,5

>=-10

>=+800

>=+1000

>=+700

 $\pi_2 = \{ (s1, move(l1, l4)), \}$ 

(s3, **move(l3,l4)**),

(s5, **move(l5,l4)**)}

(s2, wait),

(s4, wait),

### **Second Iteration**

# Policy Iteration 7: Expected Utilities for $\pi_2$

- Calculate <u>true</u> expected utilities for the <u>new</u> policy π<sub>2</sub>
  - $E(\pi 2,s1) = R(s1, move(l1,l4)) + \gamma \dots = -1 + 0.9 (0.5E(\pi 2,s1) + 0.5E(\pi 2,s4))$
  - $E(\pi 2,s2) = R(s2, wait) + \gamma E(\pi 2,s2) = -1 + 0.9 E(\pi 2,s2)$
  - $E(\pi 2, s3) = \frac{R(s3, move(13, 14))}{P} + \gamma E(\pi 2, s4) = -100 + 0.9 E(\pi 2, s4)$
  - $E(\pi 2,s4) = \frac{R(s4, wait)}{P} + \gamma E(\pi 2,s4) = +100 + 0.9 E(\pi 2,s4)$
  - $E(\pi 2,s5) = \frac{R(s5, move(15,14))}{P} + \gamma E(\pi 2,s4) = -200 + 0.9 E(\pi 2,s4)$
  - Equations to solve:
    - $0.1E(\pi 2,s2) = -1$
    - $0.1E(\pi 2, s4) = +100$
    - $E(\pi 2,s3) = -100 + 0.9E(\pi 2,s4) = -100 + 0.9*1000 = +800$
    - $E(\pi 2,s5) = -200 + 0.9E(\pi 2,s4) = -200 + 0.9*1000 = +700$
    - $E(\pi 2,s1) = -1 + 0.45 * E(\pi 2,s1) + 0.45 * E(\pi 2,s4) \rightarrow 0.55 E(\pi 2,s1) = -1 + 0.45 * E(\pi 2,s4) \rightarrow 0.55 E(\pi 2,s1) = -1 + 450 \rightarrow 0.55 E(\pi 2,s1) = +449 \rightarrow E(\pi 2,s1) = +816,3636...$

- →  $E(\pi 2, s2) = -10$
- → E(π2,s4) = +1000
- → E(π2,s3) = +800
- →  $E(\pi 2, s5) = +700$

→ E(π2,s1) = +816,36

```
\pi_2 = \{(s1, move(l1,l4), (s2, wait), (s3, move(l3,l4)), (s4, wait), (s5, move(l5,l4))\}
```

### **Policy Iteration 8: Second Policy**



### Now we have the <u>true</u> expected utilities of the second policy...

$\pi_1 = \{(s1, wait),$	E(π1,s1)=—10	$\pi_2 = \{ (s1, move(l1, l4), $	>=+ <u>444,5</u>	E(π2,s1)=+ <u>816,36</u>
(s2, wait),	$E(\pi 1, s2) = -10$	(s2, wait),	>=-10	E(π2,s2)=- 10
(s3, wait),	$E(\pi 1, s3) = -10$	(s3, move(l3,l4)),	>=+800	$E(\pi 2, s3) = +800$
(s4, wait),	$E(\pi 1, s4) = +1000$	(s4, wait),	>=+1000	$E(\pi 2, s4) = +1000$
(s5, wait)}	$E(\pi 1, s5) = -1000$	(s5, move(l5,l4))}	>=+700	$E(\pi 2, s5) = +700$

S5 wasn't so bad after all, since you can reach s4 in a single step!

S1 / s3 are even better.

S2 seems much worse in comparison, since the benefits of s4 haven't "propagated" that far.



# Policy Iteration 9: Update 2a



What is the bestElocal modificationEaccording to theEexpected utilitiesEof the current policy?E	$(\pi 2,s1) = +816,36$ $(\pi 2,s2) = -10$ $(\pi 2,s3) = +800$ $(\pi 2,s4) = +1000$ $(\pi 2,s5) = +700$		r=-101 0.8 r=-1 10.8 r=-1 101 r=-1 101 r=-101 r=-101 r=-101		
For every state s:		r=-1 sl 0.5	=99 s4 r=0		
• Let $\pi_3(s) = \operatorname{argmax}_{a \in A} Q(\pi_2, s, a)$					
• That is, find the action a that maximizes $R(s, a) + \gamma \sum_{s' \in S} P(s, a, s') E(\pi_2, s')$					
<ul><li>s1: wait</li></ul>	<b>-1</b> + 0.9 * 81	.6,36	= +733,72		
move(l1,l2)	<mark>-100</mark> +0.9*-1	0	= -109		
move(l1,l4)	<mark>—1</mark> + 0.9 * (.5	5*1000+.5*816.36)	= +816,36		
Seems best – chosen!					
<ul><li>s2: wait</li></ul>	<u> </u>	.0	= -10		
move(l2,l1)	<mark>-100</mark> + 0.9 * 81	6,36	= +634,72		
move(l2,l3)	<mark>-1</mark> + 0.9 * (0.	8*800 + 0.2*700)	= +701		

**Now** we will change the action taken at s2,

since the expected utilities for possible "next" states s1, s3, s5... have increased

# Policy Iteration 10: Update 2b



What is the best <u>local</u> modification according to the <u>expected utilities</u> of the <u>current</u> policy?	$E(\pi 2,s1) = +816,36$ $E(\pi 2,s2) = -10$ $E(\pi 2,s3) = +800$ $E(\pi 2,s4) = +1000$ $E(\pi 2,s5) = +700$		r=−1 0.8 r=−1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
For every state s:			r=99
• Let $\pi_3(s) = \arg(s)$	$\mathbf{x}_{a \in A} Q(\pi_2, s, a)$	r=-1 0.5	r=-1 r=100
That is, find the a	ction <i>a</i> that maxir	nize <mark>s R(s, a)</mark> + γ∑	$\mathbb{L}_{s' \in S} P(s, a, s') E(\pi_2, s')$
<ul> <li>s3: wait</li> </ul>	<mark>-1</mark> + 0.9 * 80	0	= +719
move(13,12)	<mark>-1</mark> + 0.9 * -1	0	= -10
move(13,14)	<mark>-100</mark> + 0.9 * 10	00	= +800
<ul> <li>s4: wait</li> </ul>	+100+0.9*10	00	= +1000
move(l4,l1)	+9 <mark>9</mark> + 0.9 * 81	6,36	= +833,72
•••			
<ul><li>s5: wait</li></ul>	<mark>-100</mark> + 0.9 * 70	0	= +530
move(15,12)	<mark>-101</mark> + 0.9 * -1	0	= -110
move(15,14)	<mark>-200</mark> + 0.9 <b>* -1</b>	000	= +700
# **Policy Iteration 11: Third Policy**

- This results in a **new policy**  $\pi_3$ 
  - True expected utilities are updated by solving an equation system
  - The algorithm will iterate once more
  - No changes will be made to the policy
  - → Termination with optimal policy!







### **Policy Iteration Algorithm**

# **Policy Iteration Algorithm**



- Policy iteration is a way to find an optimal policy π<sup>\*</sup>
  - Start with an **arbitrary** initial policy  $\pi_1$ . Then, for i = 1, 2, ...
    - Compute expected utilities  $E(\pi_i,s)$  for all s by **solving a system of equations**

Find utilities according to current policy	<ul> <li>System: For all s, E(π<sub>i</sub>,s) = Q(π<sub>i</sub>, s, π<sub>i</sub>(s)) = R(s, π<sub>i</sub>(s)) + γ∑<sub>s'∈S</sub> P(s, π<sub>i</sub>(s), s') E(π<sub>i</sub>,s')</li> <li>Result: The <u>expected utilities</u> of the "current" policy in <u>every</u> state s</li> </ul>
	Not a simple recursive calculation – the state graph is generally cyclic!
1.1.1.1	Compute an improved policy $\pi_{i+1}$ "locally" for every s
Find best local improvements	<ul> <li>π<sub>i+1</sub>(s) := argmax<sub>a ∈ A</sub> Q(π<sub>i</sub>, s, a) = argmax<sub>a ∈ A</sub> R(s, a) + γ∑<sub>s' ∈ S</sub> P(s, a, s') E(π<sub>i</sub>,s')</li> <li>Best action in any given state s given expected utilities of old policy π<sub>i</sub></li> </ul>
	If $\pi_{i+1} = \pi_i$ then exit
	No local improvement possible, so the solution is optimal

- No local improvement possible, so the solution is optimal
- Otherwise
  - This is a new policy  $\pi_{i+1}$  with **<u>new</u>** expected utilities!
  - Iterate, calculate <u>those</u> utilities, ...

# **Policy Iteration Convergence**



- **Converges** to a final answer in a finite number of iterations!
  - 1. Finite states, finite actions  $\rightarrow$  finite number of candidate policies
  - 2. An iteration can never <u>return</u> to a previous policy
    - We change which action to execute in state s only if this <u>improves expected (pseudo-)utility</u> Q for s
    - This can <u>never decrease</u> the utility for other states!
    - So utilities are monotonically strictly improving "all over"
       → no circularity possible
- Actually: <u>Polynomial</u> number of iterations!
  - But polynomial in the <u>number of states</u> (huge) not the number of objects/actions
  - May take <u>many</u> iterations, and each iteration can be slow (solving equation system)

### **Alternatives**



- Methods exist for <u>reducing the search space</u>, and for <u>approximating</u> optimal solutions (see the book)
  - Value iteration
  - Linear programming
  - Real Time Dynamic Programming

• ...

# Conclusions

jonas.kvarnstrom@liu.se - 2021



#### • Example exam topics:

- **PDB heuristics**: The main ideas of patterns, how this results in a modified planning problem, why this is faster to *solve*, how the results are used, ...
  - Given a planning problem, can you apply a pattern and find the relaxed problem?
- Landmarks: The main ideas, what a landmark *is*, how to find landmarks, how to use them in a heuristic function, ...
  - Given a planning problem, can you find n unachieved fact landmarks using the means-ends analysis algorithm?
- The **concepts** of histories, utility, discount factors, ...
- What a *policy* is / how it is defined, why we use it in some types of planning, and why a classical plan is not sufficient in these cases
  - Explain **policy iteration**, and apply 1-2 steps given a small problem instance

# TDDD48 Automated Planning (1)

#### Deeper discussions about all of these topics, and...

- Formal basis for planning
  - Alternative representations of planning problems
  - Simple and complex state transition systems
- Different principles for heuristics
- Alternative search spaces
  - Partial order planning, ...
- Extended expressivity
  - Planning with non-classical goals
- Planning with **domain knowledge** 
  - Using what you know: Temporal control rules
  - Breaking down a task into smaller parts: Hierarchical Task Networks
- **Combining planners** portfolio planning, learning planning parameters, ...
- Alternative types of planning
  - Path planning
- And so on...