

Partially Observable Markov Decision Process (POMDP)

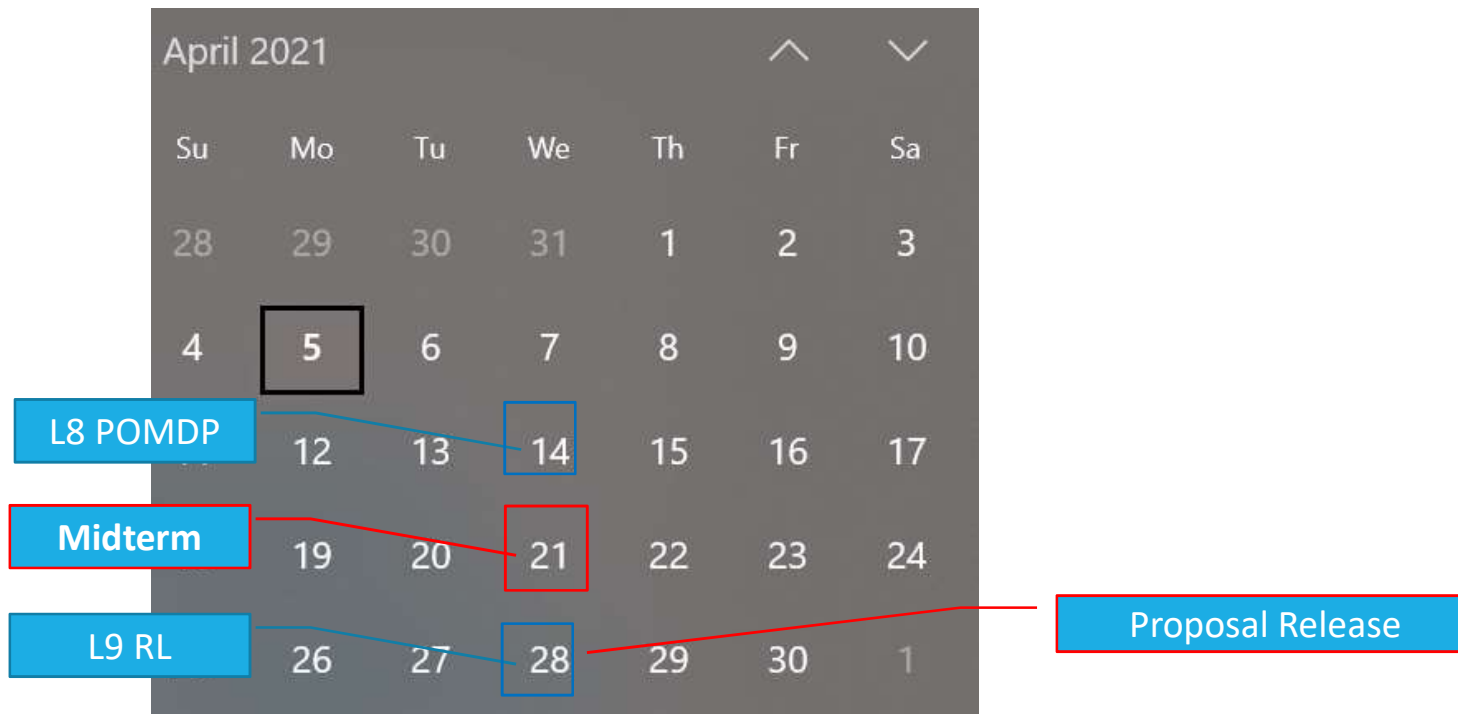
KUO-SHIH TSENG (曾國師)
DEPARTMENT OF MATHEMATICS
NATIONAL CENTRAL UNIVERSITY, TAIWAN

2021/04/14

Course Announcement

- ***Midterm (04/21/2020), 3-5pm, in M430***
 - Given a real world problem.
 - Design a perception and decision-making system for this problem using MDP , MCTS and Bayesian approaches.
- You can take one A4-size cheating sheet.
- You cannot use any electrical devices (e.g., Notebook or mobilephone), which can access to internet.
- You don't need calculators.
- You can find the **midterm_sample.pdf** on the eeclass.

Course Announcement



Course Announcement

May 2021

Su	Mo	Tu	We	Th	Fr	Sa
	3	4	5	6	7	8
	10	11	12	13	14	15
	17	18	19	20	21	22
	24	25	26	27	28	29
	31	1	2	3	4	5
	7	8	9	10	11	12
	14	15	16	17	18	19
	21	22	23	24	25	26

ROS
Tutorial

L11 NB
Perceptron

L12
Adaboost

L10 GP

L13
DL and DRL

L14
Kmeans, EM

Final Project
Presentation

Final Project
DEMO

HW3 released

Proposal DL

HW3 DL

Final Project report

Course Announcement

- In the final project, students need to propose interesting AI applications, formulate the problem, implement algorithms on a real robot (e.g., Minibot or Bebop).
 - A group should include **no more** than 2 students. If your project is large enough, you can have 3 group members (not recommended).
 - **Suggestion:** Find someone whose ability is similar to yours.
- The team members get the same credits (40%)!
 - Project Proposal (1-2 page) 10%
 - Project Presentation and Demonstration 10%
 - Project Report (4-8 pages) 20%
- You can find the sample files of final project on LMS.

Outline

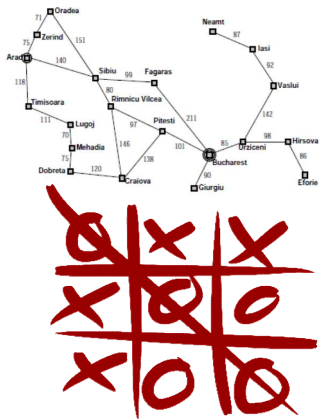
- Need
- Bayes filter
- POMDP

Outline

[Problem solving]

Search problems

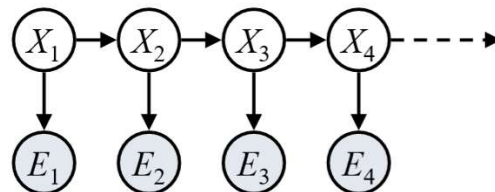
Adversarial Search



[Perception and Uncertainty]

Bayes Theorem

Bayes Filter and Smoothing

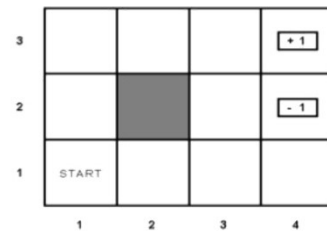
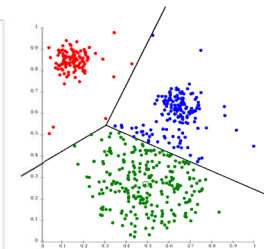
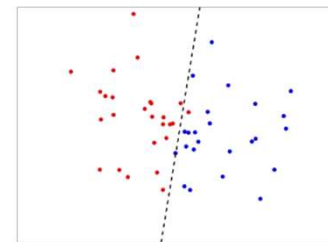


[Learning and Decision-making]

Supervised learning

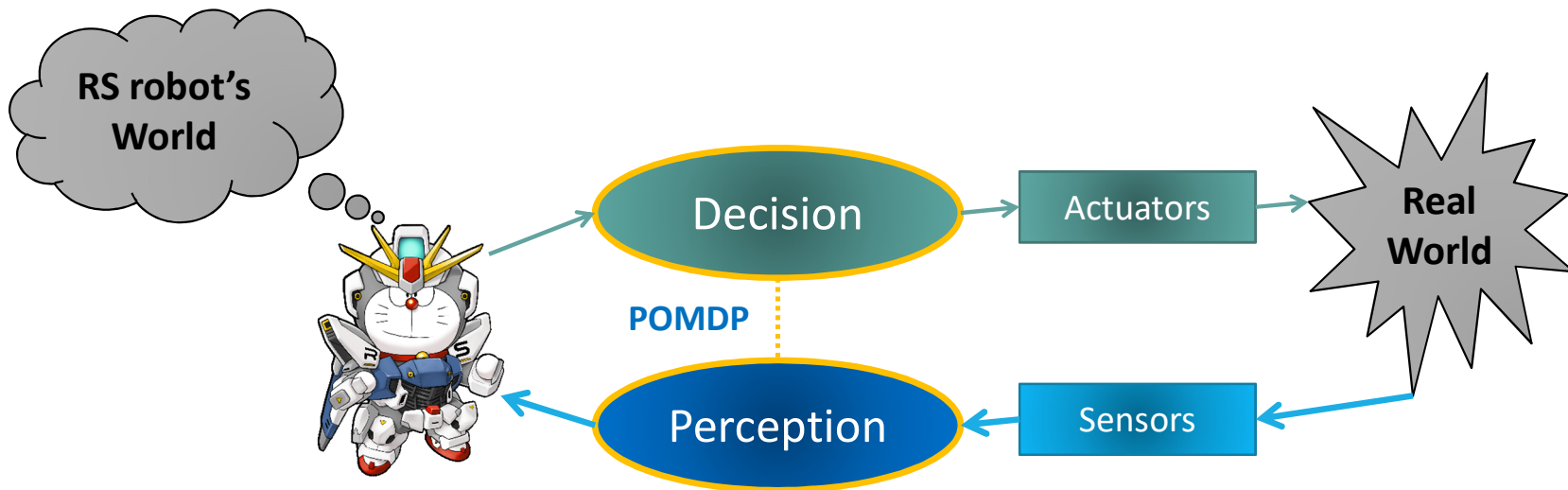
Unsupervised learning

Reinforcement learning



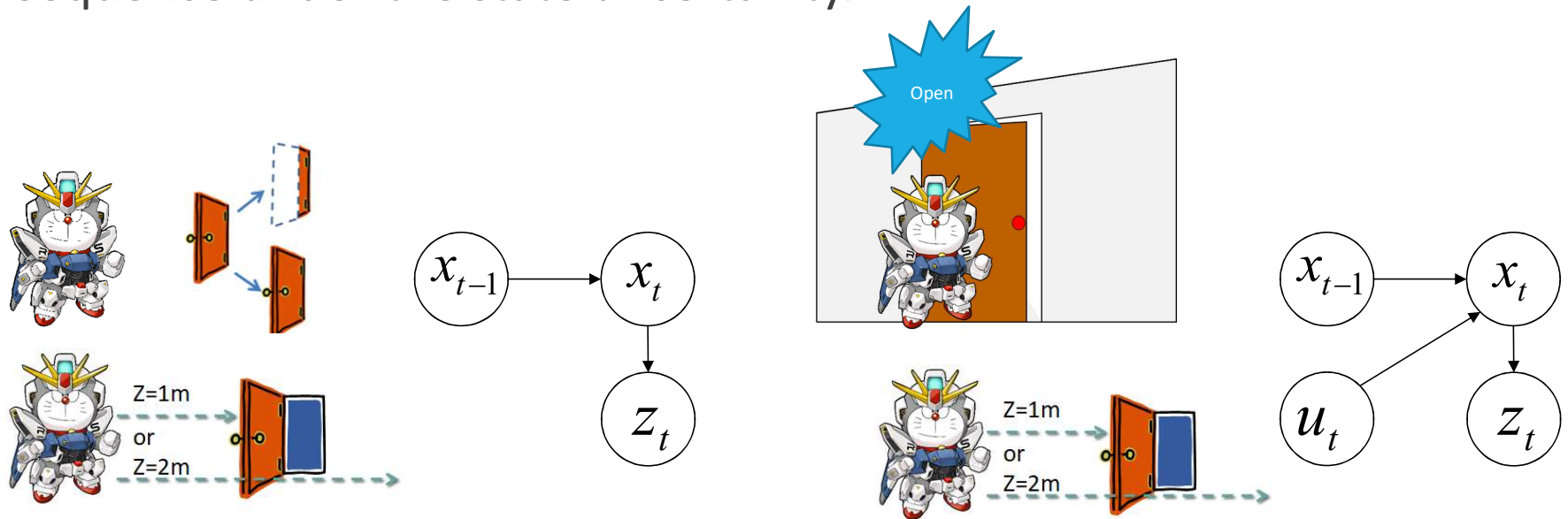
Need

- Perception
- Decision making
- Feedback



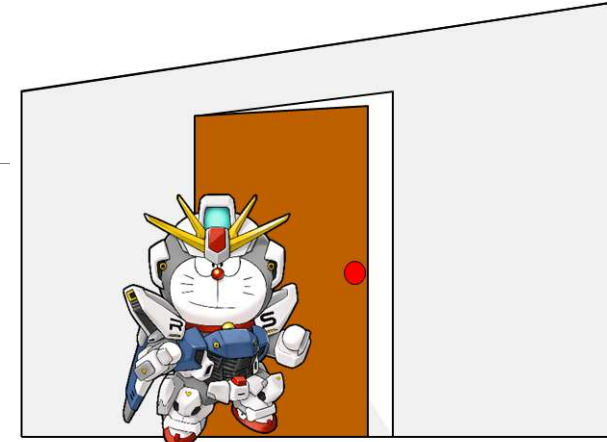
Need

- The state of MDP is deterministic.
- If the state of MDP is probabilistic, it's called **partially observable** Markov decision process (POMDP), which finds the optimal action sequence under the state uncertainty.



Need

- Example: $z=2m$, Open a door?



$$P(x_{t-1} = open) = P(x_{t-1} = close) = 0.5$$

$$P(z_t = 2 \mid x_t = open) = 0.6, P(z_t = 1 \mid x_t = open) = 0.4$$

$$P(z_t = 2 \mid x_t = close) = 0.2, P(z_t = 1 \mid x_t = close) = 0.8$$

$$P(x_t = open \mid u_t = push, x_{t-1} = open) = 1$$

$$P(x_t = close \mid u_t = push, x_{t-1} = open) = 0$$

$$P(x_t = open \mid u_t = push, x_{t-1} = close) = 0.8$$

$$P(x_t = close \mid u_t = push, x_{t-1} = close) = 0.2$$

$$P(x_t = open \mid u_t = nothing, x_{t-1} = open) = 1$$

$$P(x_t = close \mid u_t = nothing, x_{t-1} = open) = 0$$

$$P(x_t = open \mid u_t = nothing, x_{t-1} = close) = 0$$

$$P(x_t = close \mid u_t = nothing, x_{t-1} = close) = 1$$

Find : $P(x_t = open \mid u_t = nothing, z_t = 2) = ?$

$P(x_{t+1} = open \mid u_{t+1} = push, z_{t+1} = 2) = ?$

Transition matrix

Bayes filter

- Prediction

$$P(x_t | x_{t-1}, u_t, z_{t-1}) = \int \underbrace{P(x_t | x_{t-1}, u_t)}_{\text{Known (motion model)}} P(x_{t-1} | z_{t-1}) dx_{t-1}$$

- Correction

$$P(x_t | z_t) = \eta \bullet \underbrace{P(z_t | x_t)}_{\text{Known (sensor model)}} P(x_t | x_{t-1}, u_t, z_{t-1})$$

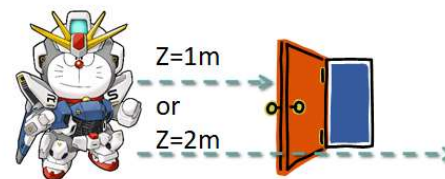
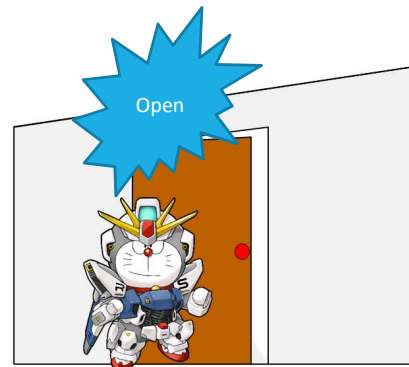
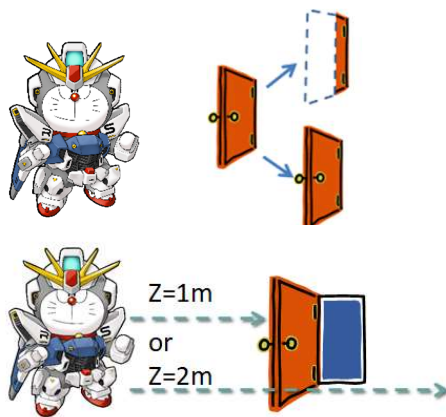
↑
constant

The robot action (u_t) and the door state (X_t) is relevant

Bayes filter

- EX:

$$P(x_2 \mid u_1 = \text{nothing}, z_1 = 2, u_2 = \text{open}, z_2 = 2) = ?$$



Bayes filter

$$P(x_{t-1} = open) = P(x_{t-1} = close) = 0.5$$

$$P(z_t = 2 \mid x_t = close) = 0.2, P(z_t = 1 \mid x_t = close) = 0.8$$

$$P(x_t = open | u_t = nothing, x_{t-1} = open) = 1$$

$$P(x_t = close | u_t = nothing, x_{t-1} = open) = 0$$

$$P(x_t = open | u_t = nothing, x_{t-1} = close) = 0$$

$$P(x_t = close | u_t = nothing, x_{t-1} = close) = 1$$

$$\begin{cases} P(x_t = o \mid x_{t-1}, z_{t-1}) = P(x_t = o \mid u_t = n, x_{t-1} = o)P(x_{t-1} = o) + P(x_t = o \mid u_t = n, x_{t-1} = c)P(x_{t-1} = c) = 0.5 \\ \quad \quad \quad \textcolor{red}{1} \qquad \qquad \qquad \textcolor{red}{0.5} \qquad \qquad \qquad \textcolor{red}{0} \qquad \qquad \qquad \textcolor{red}{0.5} \\ P(x_t = c \mid x_{t-1}, z_{t-1}) = P(x_t = c \mid u_t = n, x_{t-1} = o)P(x_{t-1} = o) + P(x_t = c \mid u_t = n, x_{t-1} = c)P(x_{t-1} = c) = 0.5 \\ \quad \quad \quad \textcolor{red}{0} \qquad \qquad \qquad \textcolor{red}{0.5} \qquad \qquad \qquad \textcolor{red}{1} \qquad \qquad \qquad \textcolor{red}{0.5} \end{cases}$$

$$P(x_t = \text{open} \mid u_t = \text{nothing}, z_t = 2) = ?$$

Bayes filter

Given:

$$P(x_{t-1} = \text{open}) = P(x_{t-1} = \text{close}) = 0.5$$

$$P(z_t = 2 \mid x_t = \text{open}) = 0.6, P(z_t = 1 \mid x_t = \text{open}) = 0.4$$

$$P(z_t = 2 \mid x_t = \text{close}) = 0.2, P(z_t = 1 \mid x_t = \text{close}) = 0.8$$

$$P(x_t = \text{open} \mid u_t = \text{push}, x_{t-1} = \text{open}) = 1$$

$$P(x_t = \text{open} \mid u_t = \text{nothing}, x_{t-1} = \text{open}) = 1$$

$$P(x_t = \text{close} \mid u_t = \text{push}, x_{t-1} = \text{open}) = 0$$

$$P(x_t = \text{close} \mid u_t = \text{nothing}, x_{t-1} = \text{open}) = 0$$

$$P(x_t = \text{open} \mid u_t = \text{push}, x_{t-1} = \text{close}) = 0.8$$

$$P(x_t = \text{open} \mid u_t = \text{nothing}, x_{t-1} = \text{close}) = 0$$

$$P(x_t = \text{close} \mid u_t = \text{push}, x_{t-1} = \text{close}) = 0.2$$

$$P(x_t = \text{close} \mid u_t = \text{nothing}, x_{t-1} = \text{close}) = 1$$

$$P(x_t \mid z_t) = \eta \bullet P(z_t \mid x_t)P(x_t \mid x_{t-1}, z_{t-1})$$

$$\left\{ \begin{array}{l} P(x_t = o \mid z_t) = \eta \bullet \underset{0.6}{P(z_t \mid x_t = o)} \underset{0.5}{P(x_t = o \mid x_{t-1}, u_t, z_{t-1})} = 0.75 \end{array} \right.$$

$$\left\{ \begin{array}{l} P(x_t = c \mid z_t) = \eta \bullet \underset{0.2}{P(z_t \mid x_t = c)} \underset{0.5}{P(x_t = c \mid x_{t-1}, u_t, z_{t-1})} = 0.25 \end{array} \right.$$

$$(0.4\eta = 1, \eta = 2.5)$$

Bayes filter

$$P(x_{t-1} = open) = P(x_{t-1} = close) = 0.5$$

$$P(z_t = 2 \mid x_t = close) = 0.2, P(z_t = 1 \mid x_t = close) = 0.8$$

$$P(x_t = open | u_t = nothing, x_{t-1} = open) = 1$$

$$P(x_t = close | u_t = nothing, x_{t-1} = open) = 0$$

$$P(x_t = open | u_t = nothing, x_{t-1} = close) = 0$$

$$P(x_t = close | u_t = nothing, x_{t-1} = close) = 1$$

$$\begin{cases} P(x_t = o \mid x_{t-1}, z_{t-1}) = P(x_t = o \mid u_t = n, x_{t-1} = o)P(x_{t-1} = o) + P(x_t = o \mid u_t = p, x_{t-1} = c)P(x_{t-1} = c) = 0.95 \\ \quad \quad \quad \textcolor{red}{1} \qquad \qquad \qquad \textcolor{red}{0.75} \qquad \qquad \qquad \textcolor{red}{0.8} \qquad \qquad \qquad \textcolor{red}{0.25} \\ P(x_t = c \mid x_{t-1}, z_{t-1}) = P(x_t = c \mid u_t = n, x_{t-1} = o)P(x_{t-1} = o) + P(x_t = c \mid u_t = p, x_{t-1} = c)P(x_{t-1} = c) = 0.05 \\ \quad \quad \quad \textcolor{red}{0} \qquad \qquad \qquad \textcolor{red}{0.75} \qquad \qquad \qquad \textcolor{red}{0.2} \qquad \qquad \qquad \textcolor{red}{0.25} \end{cases}$$

$$P(x_{t+1} = open \mid u_{t+1} = push, z_{t+1} = 2) = ?$$
$$P(x_{t-1} = open) = P(x_{t-1} = close) = 0.5$$

$$P(z_t = 2 \mid x_t = close) = 0.2, P(z_t = 1 \mid x_t = close) = 0.8$$

$$P(x_t = open \mid u_t = nothing, x_{t-1} = open) = 1$$

$$P(x_t = close | u_t = nothing, x_{t-1} = open) = 0$$

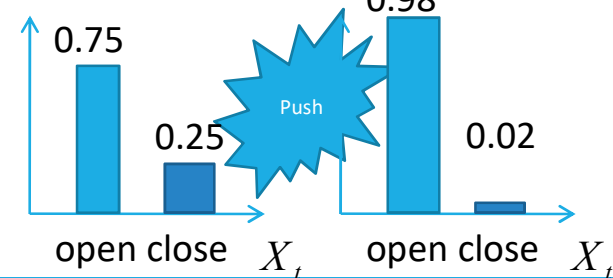
$$P(x_t = open | u_t = nothing, x_{t-1} = close) = 0$$

$$P(x_t = close | u_t = nothing, x_{t-1} = close) = 1$$

$$\overline{P(X_t \mid Z_t = 2)}$$

$$P(\overline{X}_t \mid Z_t = 2)$$

$$\left[P(x_t = c \mid z_t) = \eta \bullet \underset{0.2}{P(z_t \mid x_t = c)} \underset{0.05}{P(x_t = c \mid x_{t-1}, u_t, z_{t-1})} = 0.02 \right]$$



MDP

- MDP is a model for finding sequential optimal decisions.
 - State: fully observable
 - State transition: stochastic (Motion model)

[*GIVEN*]

S : state

A : action

$P(s'|s, a)$: Transition probability

$R(s, a)$: reward

γ : discount

[*Find*]

$\pi^* = \arg \max U^\pi(s)$

$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

π : policy

π^* : optimal policy

U : utility

POMDP

- POMDP is a model for finding sequential optimal decisions.
 - State: **stochastic** (Sensor model)
 - State transition: stochastic (Motion model)

[*GIVEN*]

S : state

A : action

$P(s'|s, a)$: Transition probability

$R(s, a)$: reward

γ : discount

z : a set of measurement

$P(z | s)$: sensor model

[*Find*]

$\pi^* = \arg \max U^\pi(s)$

$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$

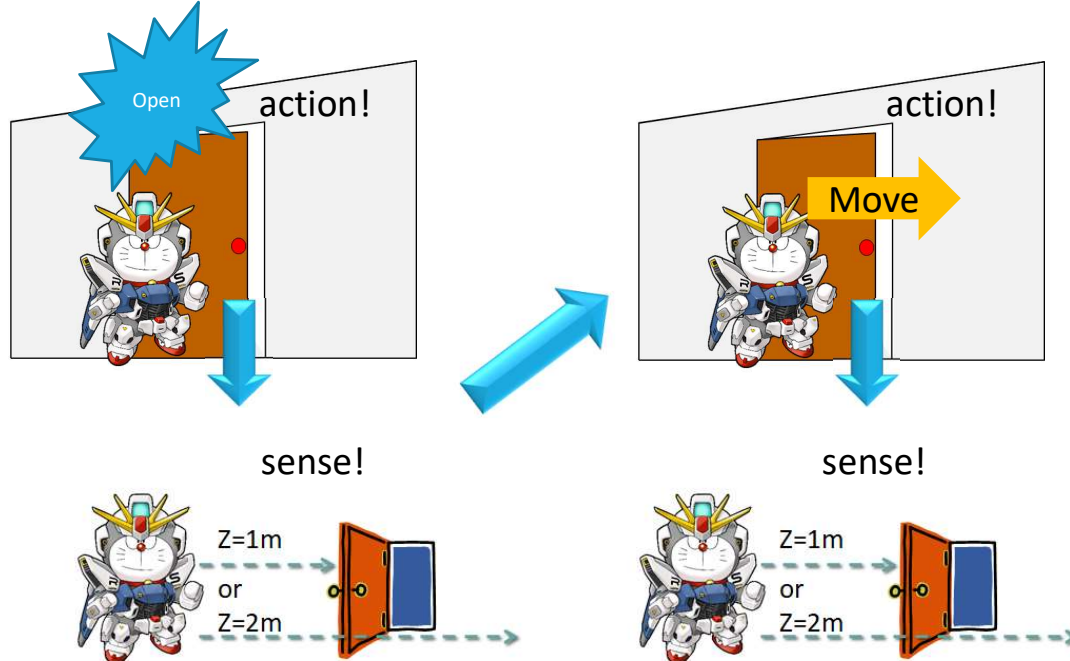
π : policy

π^* : optimal policy

U : utility

POMDP

- POMDP is a model for finding sequential optimal decisions.
 - State: **stochastic** (The robot has to *sense* the state via sensors)
 - State transition: stochastic



[Find]

$$\pi^* = \arg \max U^\pi(s)$$

$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

POMDP

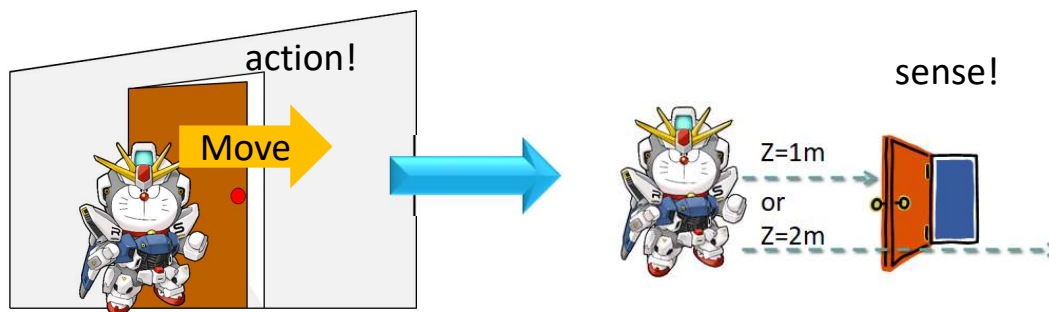
- POMDP is a model for finding sequential optimal decisions.
 - State: **stochastic** (The robot has to *sense* the state via sensors)
 - State transition: stochastic

3 steps of POMDP :

1. Given current belief state b , execute the action $a = \pi^*(b)$.
2. Receive measurement z .
3. Set the current belief state to $Foward(b, a, z)$ and repeat.

Decision

Perception



POMDP

- Let's look at an example of POMDP

[GIVEN]

$S \in \{Close, Open\}$

$A \in \{Move, Stay, Open\}$

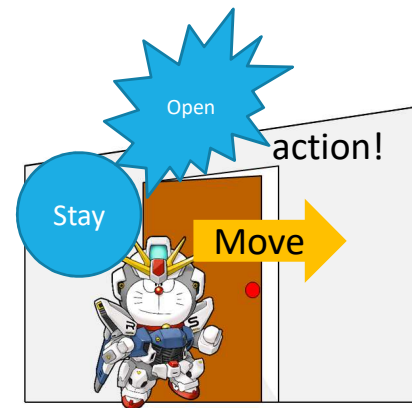
$P(s'|s, a)$: Transition probability

$R(s, a)$: reward

γ : discount

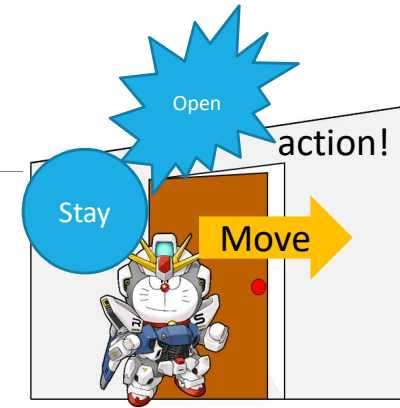
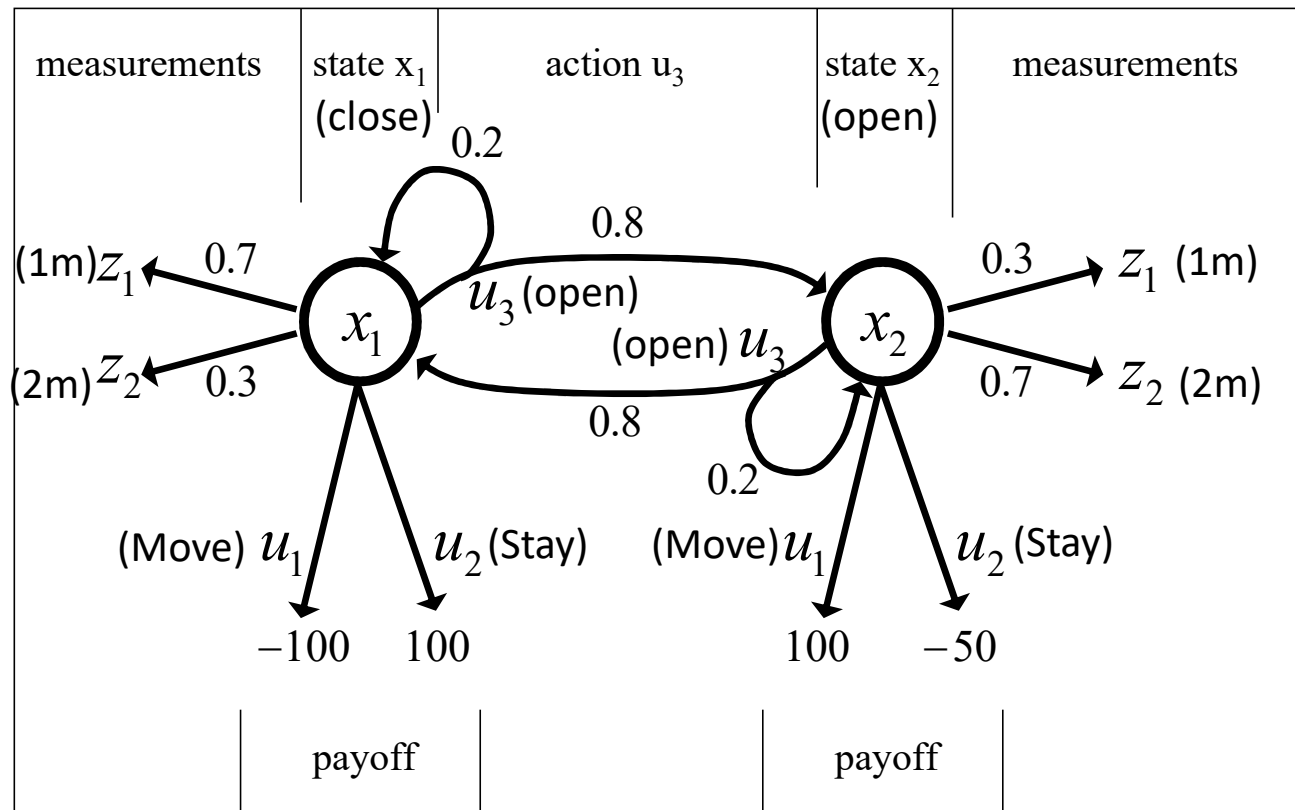
z : a set of measurement

$P(z | s)$: sensor model



Let's check the state machine of this problem!

POMDP



Example from:
Probabilistic Robotics

POMDP

- The actions u_1 and u_2 are terminal actions. The action u_3 is a sensing action that potentially leads to a state transition.
- The horizon is finite and $\gamma=1$.

Motion Model:

$$P(x_1 | x_1, u_3) = 0.2, P(x_2 | x_1, u_3) = 0.8$$

$$P(x_1 | x_2, u_3) = 0.8, P(x_2 | x_2, u_3) = 0.2$$

Sensor Model:

$$P(z_1 | x_1) = 0.7, P(z_2 | x_1) = 0.3$$

$$P(z_1 | x_2) = 0.3, P(z_2 | x_2) = 0.7$$

Reward:

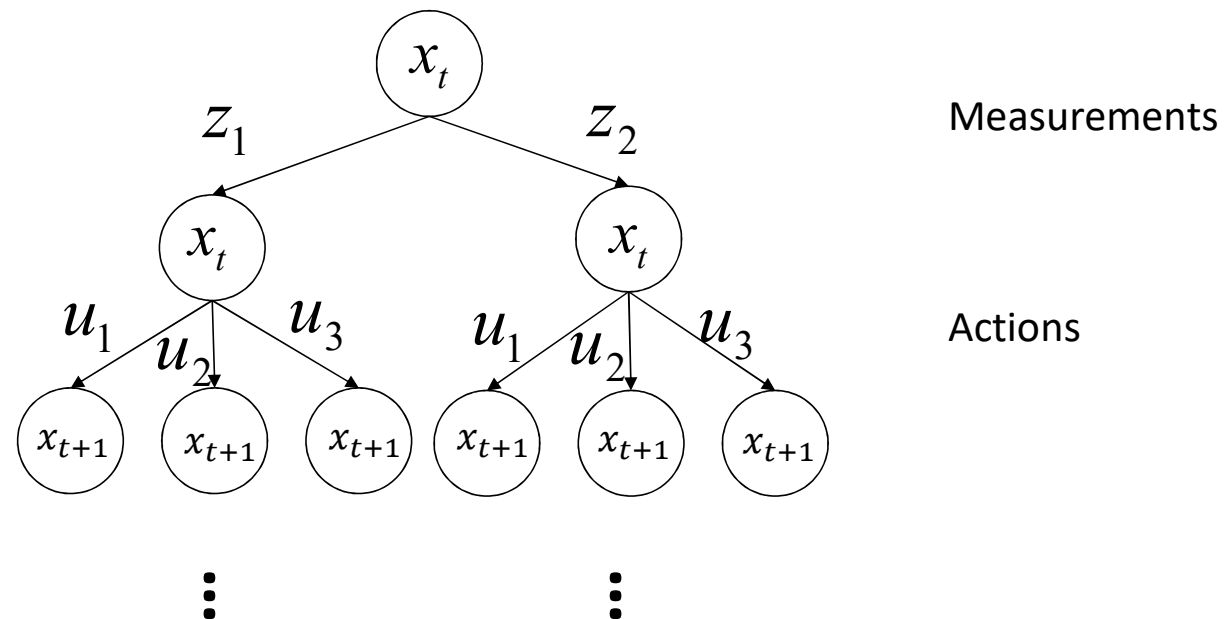
$$r(x_1, u_1) = -100, r(x_2, u_1) = +100$$

$$r(x_1, u_2) = +100, r(x_2, u_2) = -50$$

$$r(x_1, u_3) = -1, r(x_2, u_3) = -1$$

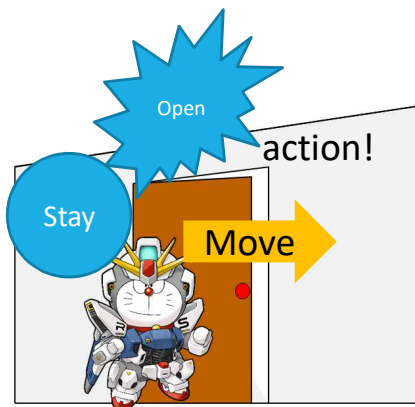
POMDP

- Search tree of POMDP

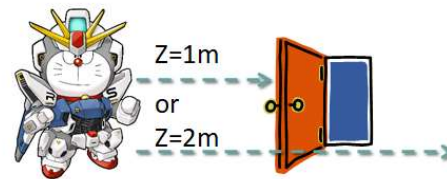


POMDP

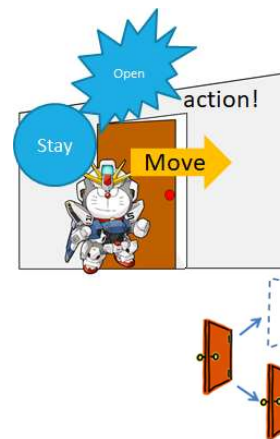
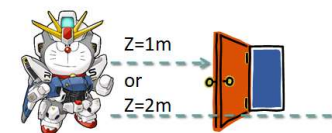
- **Action only**



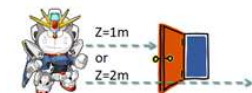
- Sense
- Action



- Sense
- Action
- Transition



- Sense
- Action
- Transition
- H horizons



• • • • •

POMDP

- In MDPs, the payoff (or return) depended on the state of the system.
- In POMDPs, however, the true state is not exactly known.
- Therefore, we compute the **expected payoff** by **integrating/sum over all states**:

$$\begin{aligned} r(b, u) &= E_x[r(x, u)] \\ &= \sum_x r(x, u)P(x) \\ &= P(x_1)r(x_1, u) + P(x_2)r(x_2, u) \end{aligned}$$

POMDP

- If we are totally certain that we are in state x_1 and execute action u_1 , we receive a reward of -100
- If, on the other hand, we definitely know that we are in x_2 and execute u_1 , the reward is +100.
- In between it is the linear combination of the extreme values weighted by the probabilities

$$\begin{aligned} r(b, u_1) &= -100P(x_1) + 100P(x_2) \\ &= -100p_1 + 100(1 - p_1) \end{aligned}$$

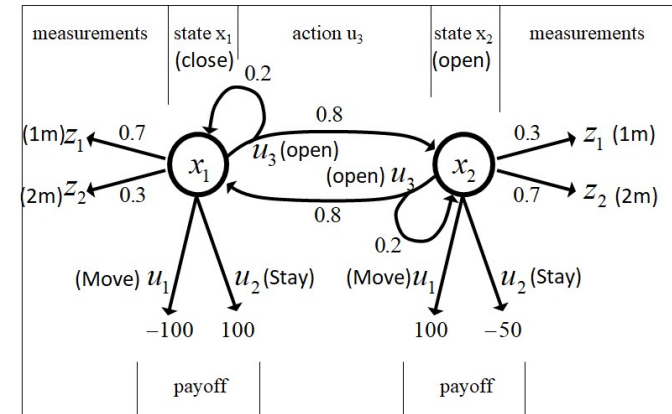
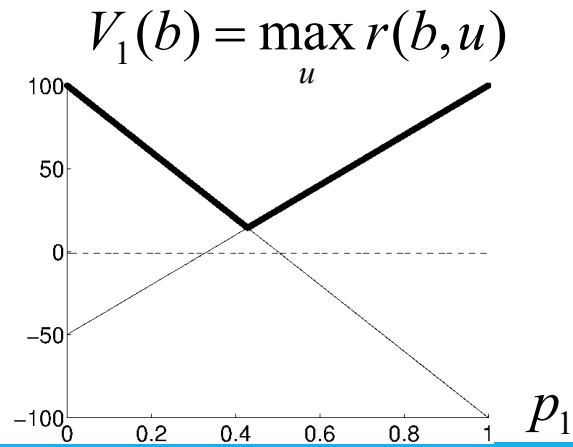
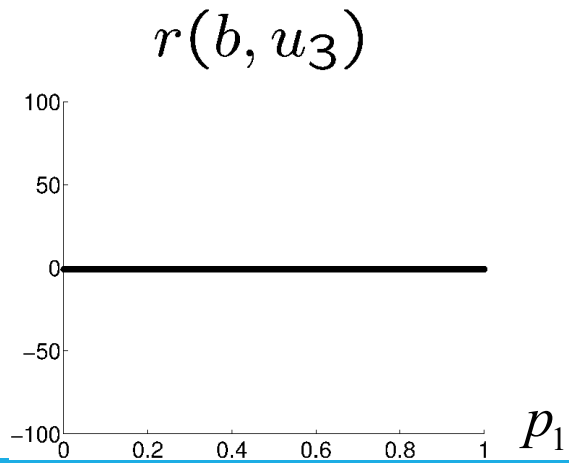
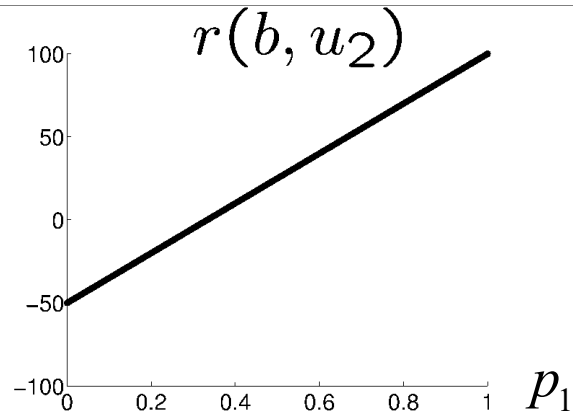
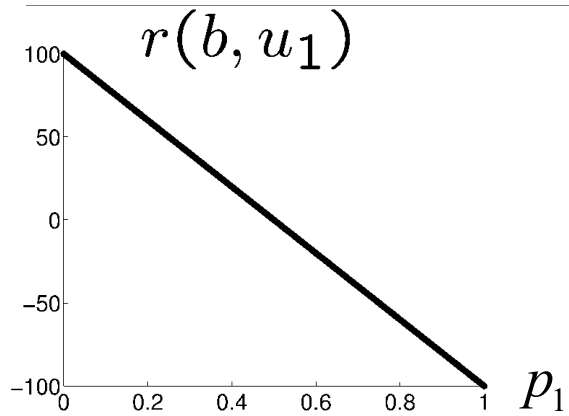
$$r(b, u_2) = +100p_1 - 50(1 - p_1)$$

$$r(b, u_3) = -1$$

$$P(x_1) = p_1$$

Let b be the belief of the agent about the state under consideration

POMDP



POMDP

- The Resulting Policy for $T=1$
- Given we have a finite POMDP with $T=1$, we would use $V_I(b)$ to determine the optimal policy.
- In our example, the optimal policy for $T=1$ is

$$\pi_1(b) = \begin{cases} u_1 & \text{if } p_1 \leq \frac{3}{7} \\ u_2 & \text{if } p_1 > \frac{3}{7} \end{cases}$$

- This is the upper thick graph in the diagram.

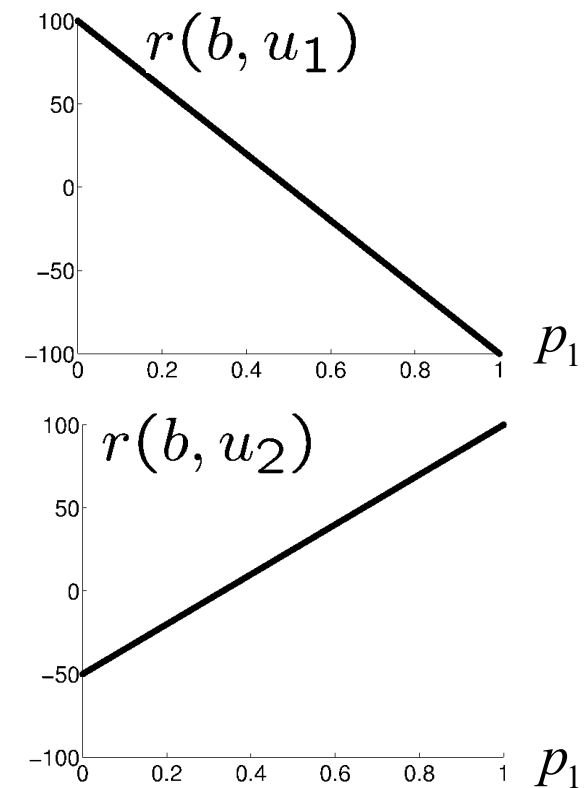
POMDP

- The resulting value function $V_I(b)$ is the maximum of the three functions at each point

$$V_1(b) = \max_u r(b, u)$$

$$= \max \begin{cases} -100p_1 + 100(1-p_1) & \leftarrow \text{U1} \\ 100p_1 - 50(1-p_1) & \leftarrow \text{U2} \\ -1 & \leftarrow \text{U3} \end{cases}$$

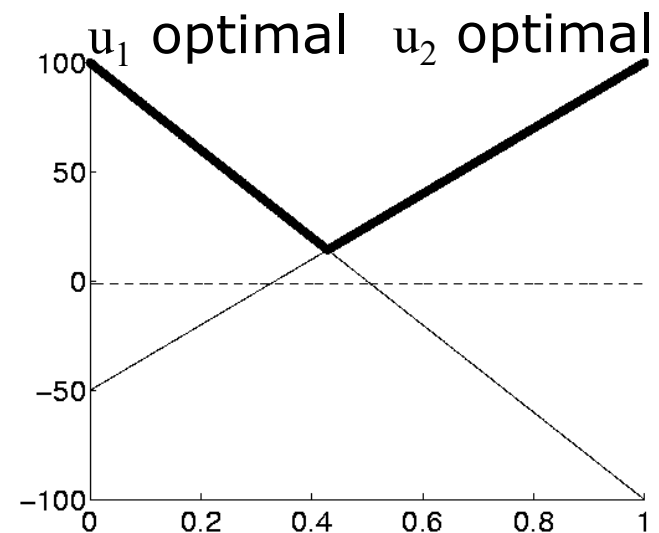
- It is piecewise linear and convex.



POMDP

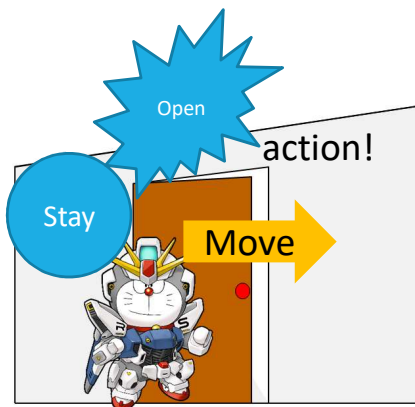
- Pruning:
- If we carefully consider $V_I(b)$, we see that only the first two components contribute.
- The third component can therefore safely be pruned away from $V_I(b)$.

$$V_I(b) = \max \begin{cases} -100p_1 + 100(1-p_1) \\ 100p_1 - 50(1-p_1) \end{cases}$$

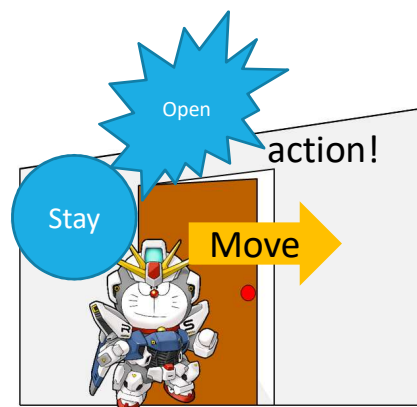
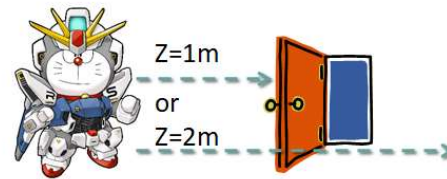


POMDP

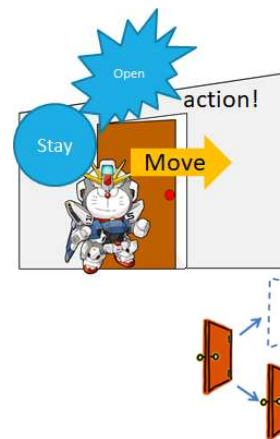
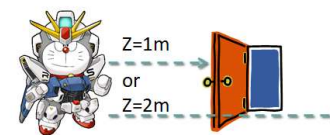
- Action only



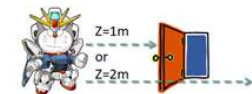
- **Sense**
- **Action**



- Sense
- Action
- Transition



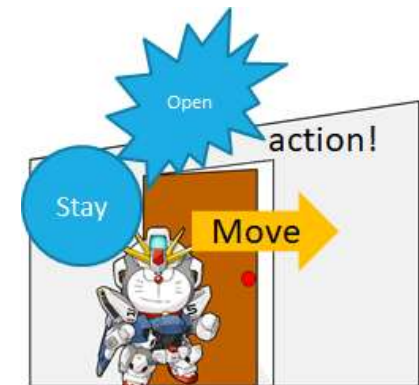
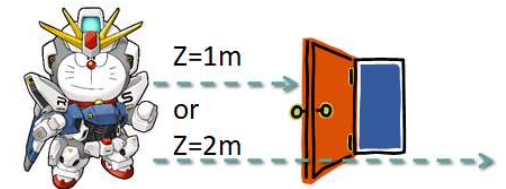
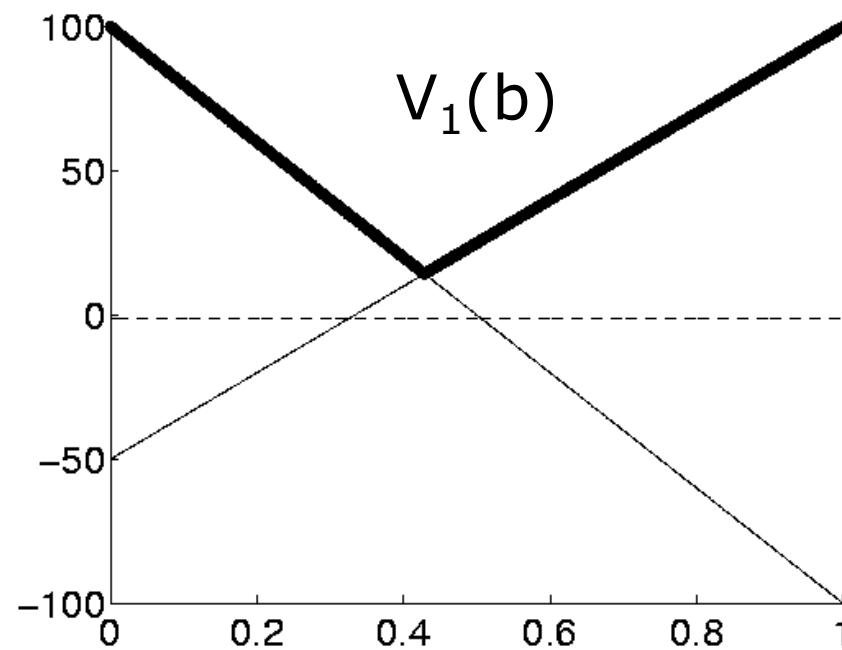
- Sense
- Action
- Transition
- H horizons



• • • • •

POMDP

- Assume the robot can make an observation before deciding on an action. → Bayes theorem!



POMDP

- Assume the robot can make an observation before deciding on an action.
- Suppose the robot perceives z_1 for which $p(z_1 | x_1)=0.7$ and $p(z_1 | x_2)=0.3$.
- Given the observation z_1 we update the belief using Bayes rule.

$$p'_1 = \frac{0.7 p_1}{p(z_1)}$$

$$p'_2 = \frac{0.3(1 - p_1)}{p(z_1)}$$

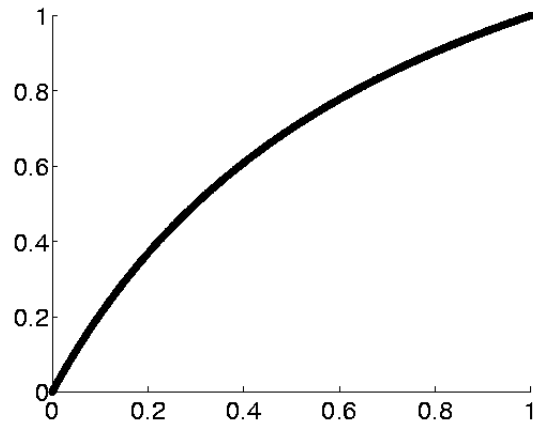
$$p(z_1) = 0.7 p_1 + 0.3(1 - p_1) = 0.4 p_1 + 0.3$$

$$P(x | z) = \frac{P(z | x)P(x)}{P(z)}$$

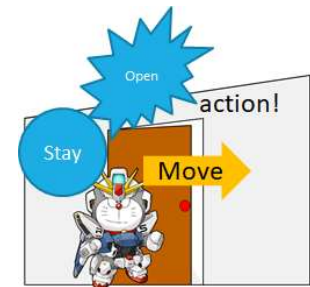
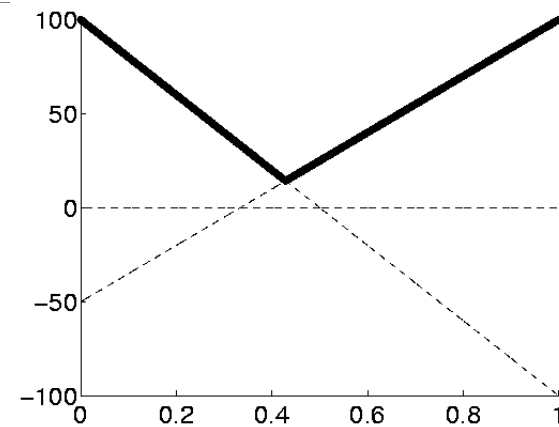
POMDP

- Value function

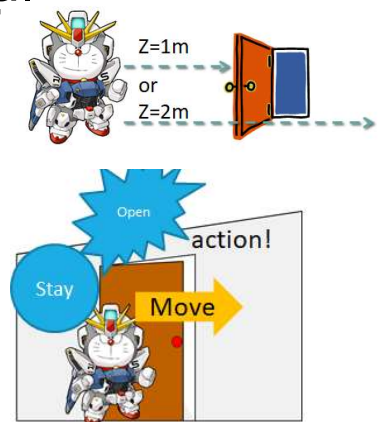
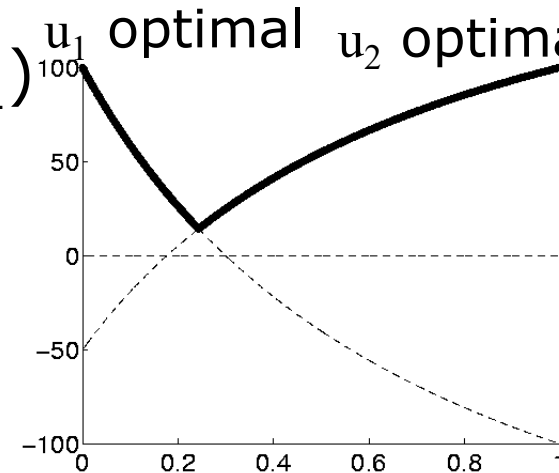
$$r(b|z_1, u_2)$$



$$V_1(b)$$



$$V_1(b|z_1) \quad u_1 \text{ optimal} \quad u_2 \text{ optimal}$$

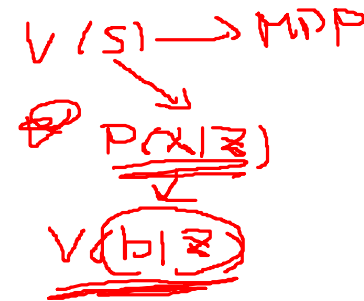


POMDP

- Assume the robot can make an observation before deciding on an action.
- **Suppose the robot perceives z_1** for which $p(z_1 | x_1)=0.7$ and $p(z_1 | x_2)=0.3$.
- Given the observation z_1 we update the belief using Bayes rule.
- Thus $V_1(b | z_1)$ is given by

$$V_1(b | z_1) = \max \begin{cases} -100 \frac{0.7 p_1}{P(z_1)} + 100 \frac{0.3(1-p_1)}{P(z_1)} \\ 100 \frac{0.7 p_1}{P(z_1)} - 50 \frac{0.3(1-p_1)}{P(z_1)} \end{cases}$$

$$= \frac{1}{P(z_1)} \max \begin{cases} 70 p_1 + 30(1-p_1) \\ 70 p_1 - 15(1-p_1) \end{cases}$$



POMDP

- Expected value after measuring:
- Since we do not know in advance what the *next* measurement will be, we have to compute the expected belief

$$\begin{aligned}\bar{V}_1(b) &= E_z[V_1(b | z)] = \sum_{i=1}^2 p(z_i) V_1(b | z_i) \\ &= \sum_{i=1}^2 p(z_i) V_1\left(\frac{p(z_i | x_1) p_1}{p(z_i)}\right) \\ &= \sum_{i=1}^2 V_1(p(z_i | x_1) p_1)\end{aligned}$$

POMDP

- Since we do not know in advance what the next measurement will be, we have to compute the expected belief

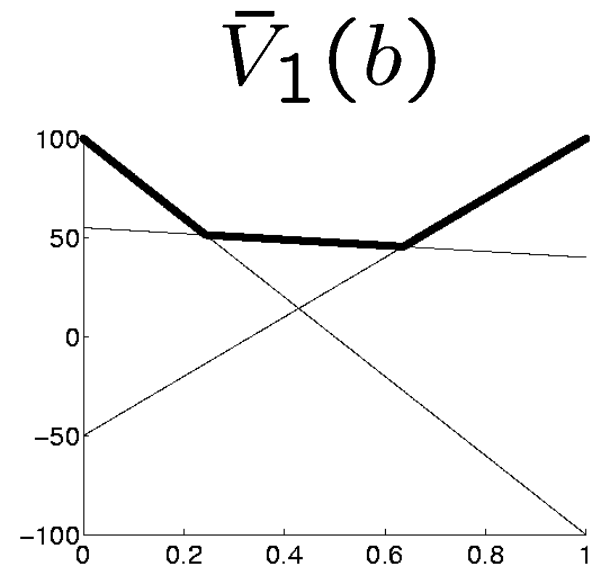
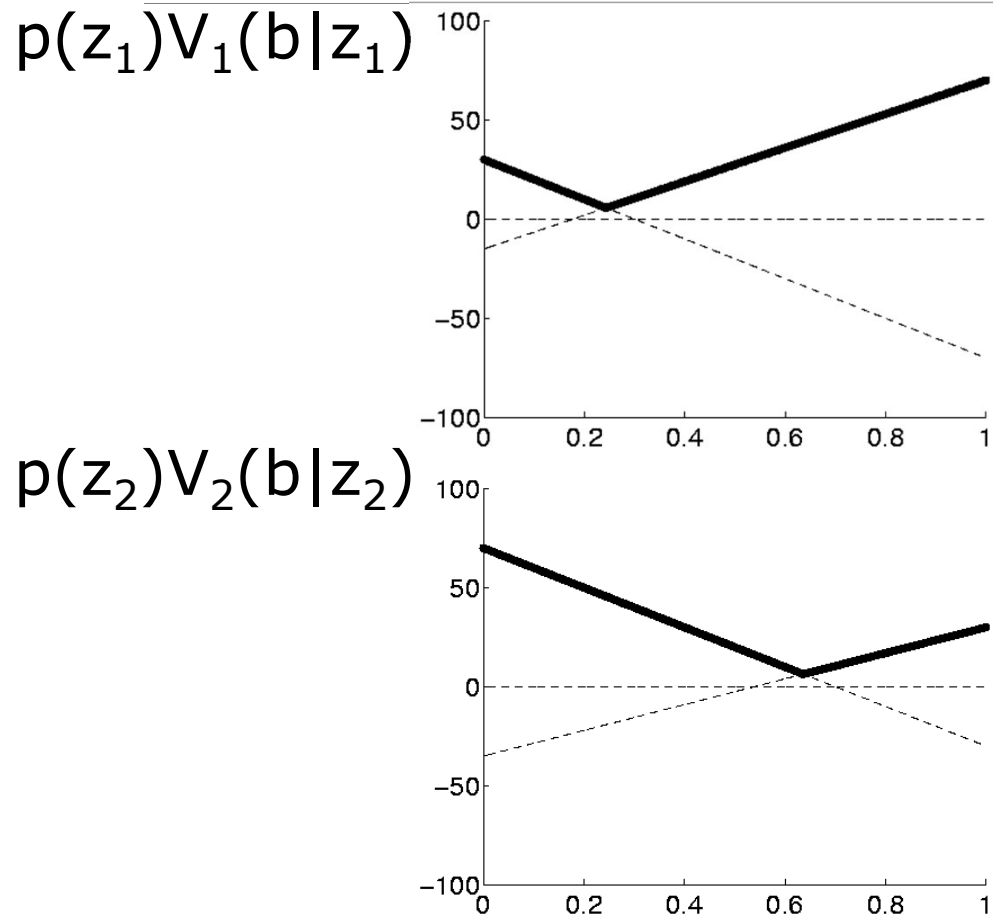
$$\begin{aligned}\bar{V}_1(b) &= E_z[V_1(b | z)] = \sum_{i=1}^2 p(z_i) V_1(b | z_i) \\ &= \max \left\{ \begin{array}{l} -70p_1 + 30(1-p_1) \\ 70p_1 - 15(1-p_1) \end{array} \right\} \\ &\quad + \max \left\{ \begin{array}{l} -30p_1 + 70(1-p_1) \\ 30p_1 - 35(1-p_1) \end{array} \right\}\end{aligned}$$

POMDP

- The four possible combinations yield the following function which then can be simplified and pruned.

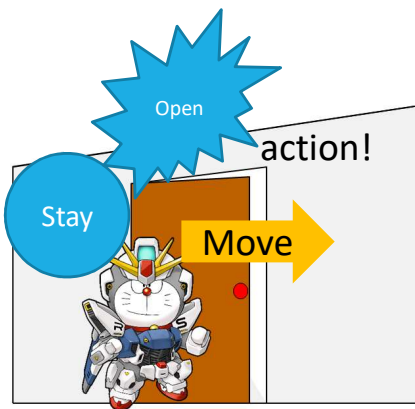
$$\begin{aligned}\bar{V}_1(b) &= \max \left\{ \begin{array}{l} -70p_1 + 30(1-p_1) - 30p_1 + 70(1-p_1) \\ -70p_1 + 30(1-p_1) + 30p_1 - 35(1-p_1) \\ 70p_1 - 15(1-p_1) - 30p_1 + 70(1-p_1) \\ 70p_1 - 15(1-p_1) + 30p_1 - 35(1-p_1) \end{array} \right\} \\ &= \max \left\{ \begin{array}{l} -100p_1 + 100(1-p_1) \\ +40p_1 + 55(1-p_1) \\ +100p_1 - 50(1-p_1) \end{array} \right\}\end{aligned}$$

POMDP

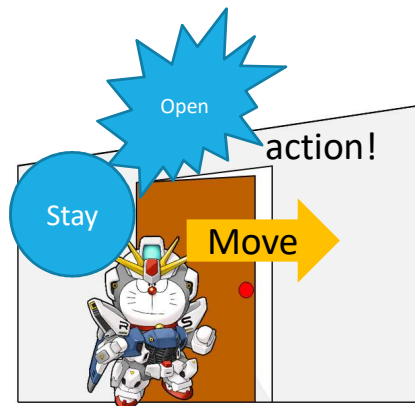
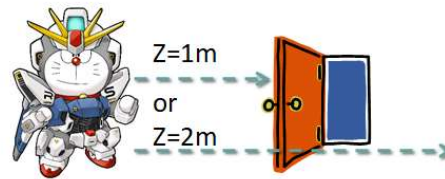


POMDP

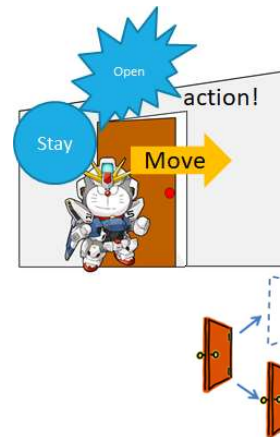
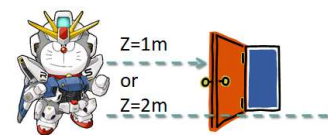
- Action only



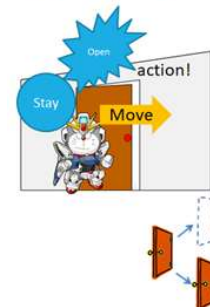
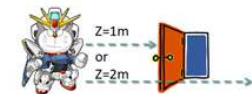
- Sense
- Action



- **Sense**
- **Action**
- **Transition**



- Sense
- Action
- Transition
- H horizons

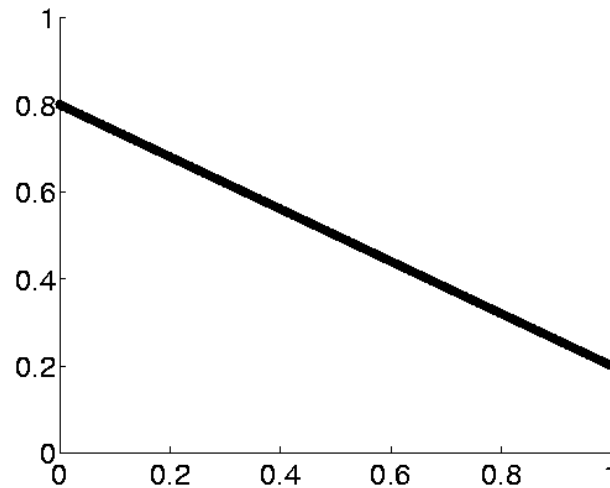


• • • • •

POMDP

- State Transitions (Prediction):
- When the agent selects u_3 its state potentially changes.
- When computing the value function, we have to take these potential state changes into account.

$$\begin{aligned} p_1' &= E_x[P(x_1 | x, u_3)] \\ &= \sum_{i=1}^2 P(x_1 | x_i, u_3) p_i \\ &= 0.2 p_1 + 0.8(1 - p_1) \\ &= 0.8 - 0.6 p_1 \end{aligned}$$



$$\begin{bmatrix} 0.2 & 0.8 \\ 0.8 & 0.2 \end{bmatrix} \begin{bmatrix} p_1 \\ 1 - p_1 \end{bmatrix}$$

POMDP

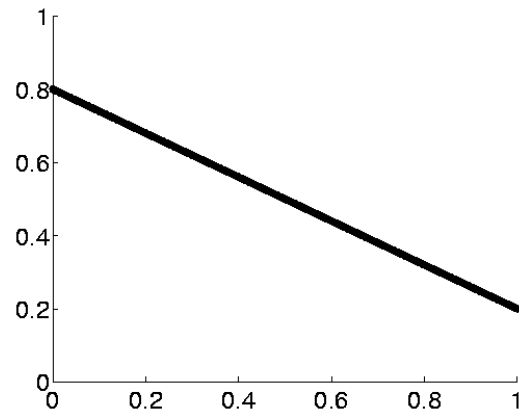
- Resulting Value Function after executing u_3
- Taking the state transitions into account, we finally obtain.

$$\bar{V}_1(b) = \max \left\{ \begin{array}{l} -100p_1 + 100(1-p_1) \\ +40p_1 + 55(1-p_1) \\ +100p_1 - 50(1-p_1) \end{array} \right\}$$

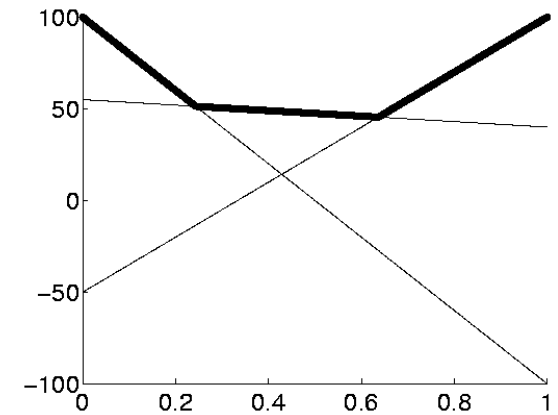
$$\bar{V}_1(b | u_3) = \max \left\{ \begin{array}{l} 60p_1 - 60(1-p_1) \\ 52p_1 + 43(1-p_1) \\ -20p_1 + 70(1-p_1) \end{array} \right\}$$

POMDP

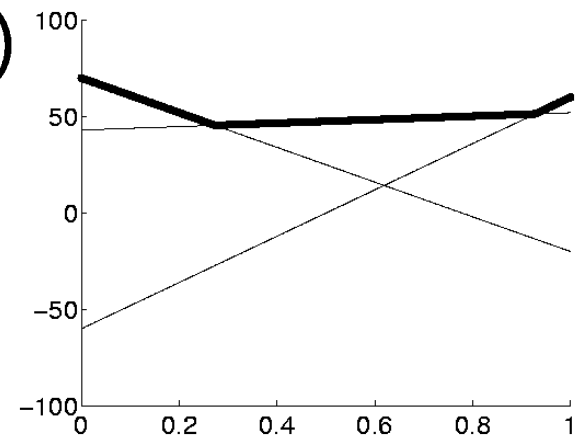
$$p_1' = 0.8 - 0.6p_1$$



$$\bar{V}_1(b)$$

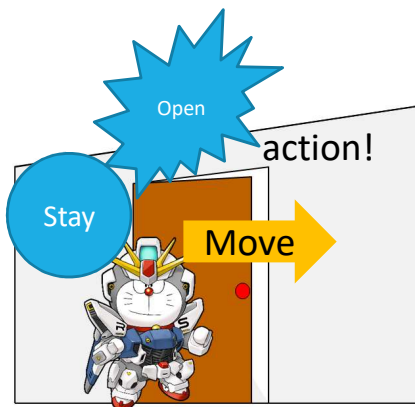


$$\bar{V}_1(b | u_3)$$

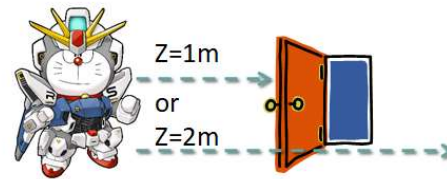


POMDP

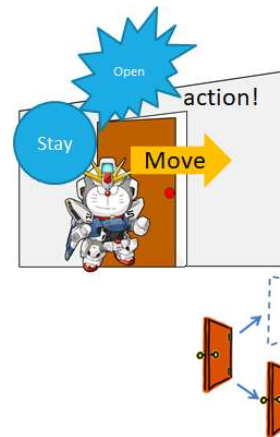
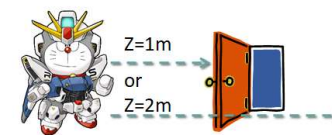
- Action only



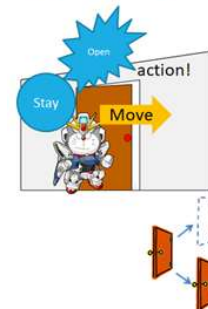
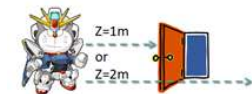
- Sense
- Action



- Sense
- Action
- Transition



- **Sense**
- **Action**
- **Transition**
- **H horizons**



• • • • •

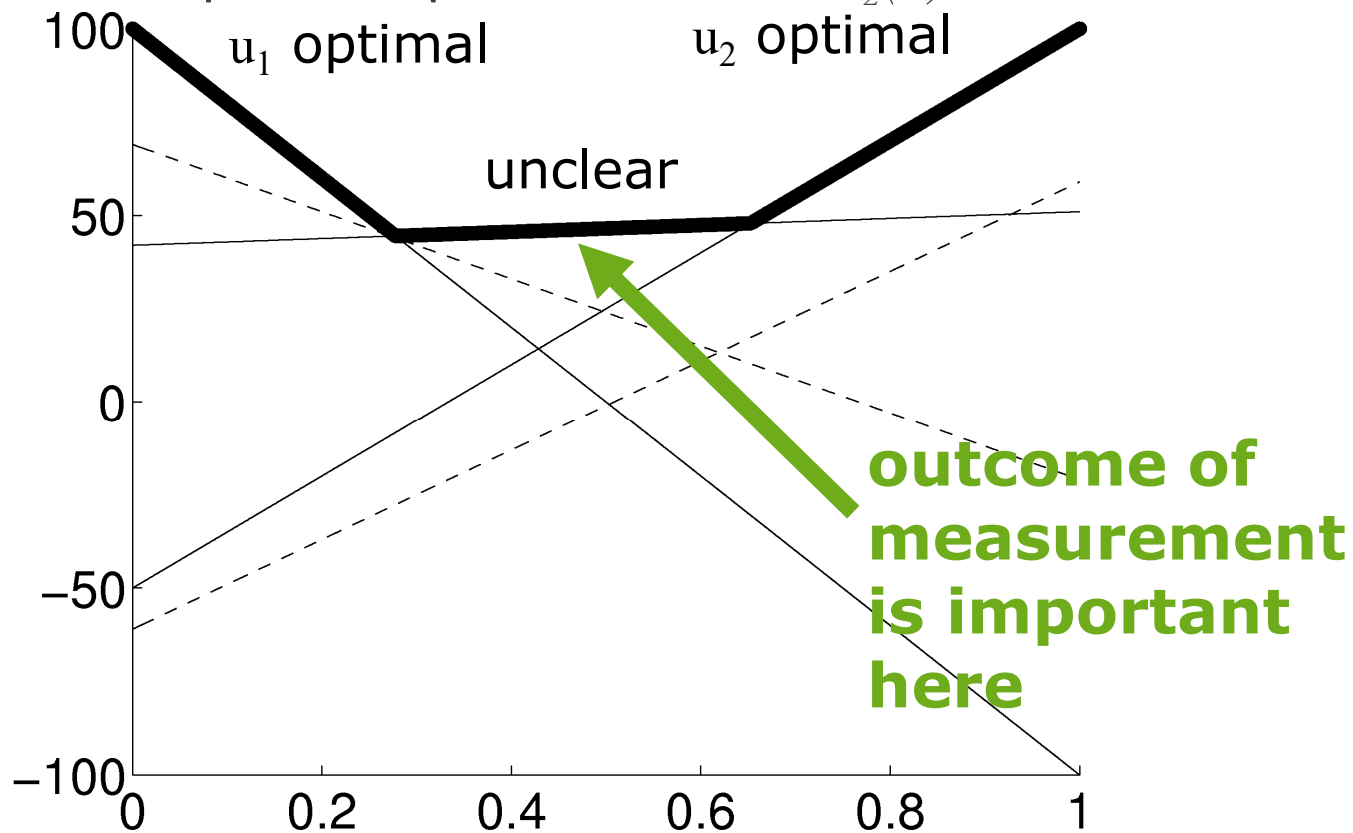
POMDP

- Value Function for $T=2$:
- Taking into account that the agent can either directly perform u_1 or u_2 or first u_3 and then u_1 or u_2 , we obtain (after pruning)

$$\bar{V}_2(b) = \max \left\{ \begin{array}{l} -100p_1 + 100(1-p_1) \\ 100p_1 - 50(1-p_1) \\ 51p_1 + 42(1-p_1) \end{array} \right\}$$

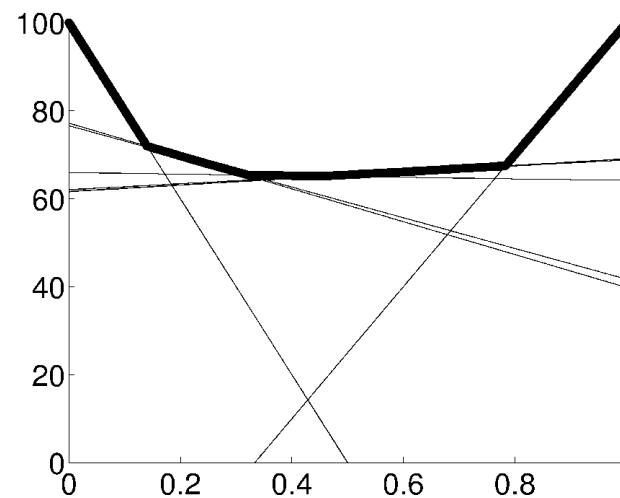
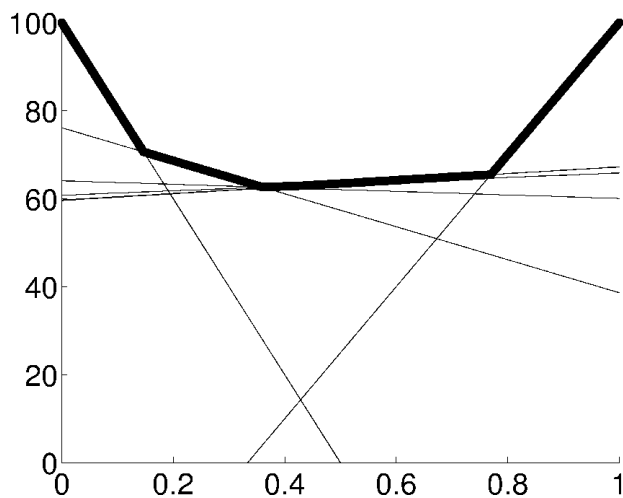
POMDP

- Graphical Representation of $V_2(b)$

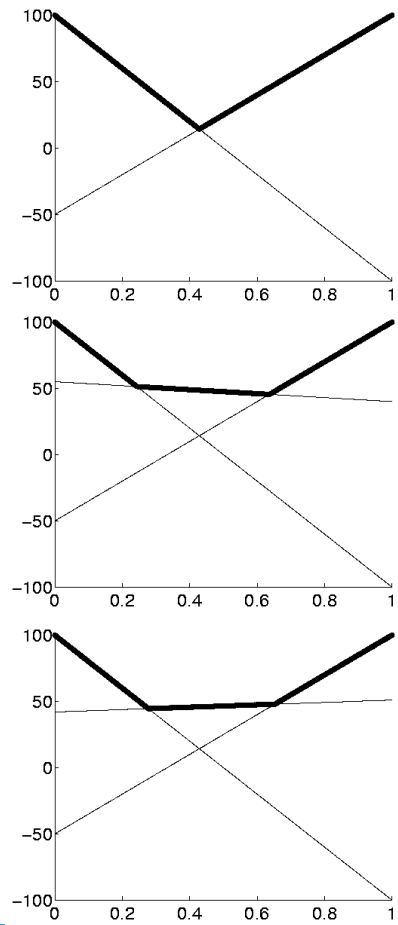


POMDP

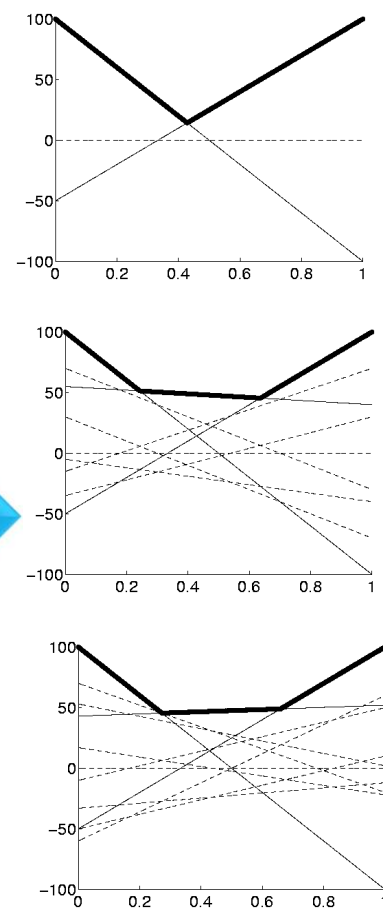
- Deep Horizons and Pruning:
- We have now completed a full backup in belief space.
- This process can be applied recursively.
- The value functions for $T=10$ and $T=20$ are



POMDP



More horizons



```

1:  Algorithm POMDP( $T$ ):
2:       $\Upsilon = (0, \dots, 0)$ 
3:      for  $\tau = 1$  to  $T$  do
4:           $\Upsilon' = \emptyset$ 
5:          for all  $(u'; v_1^k, \dots, v_N^k)$  in  $\Upsilon$  do
6:              for all control actions  $u$  do
7:                  for all measurements  $z$  do
8:                      for  $j = 1$  to  $N$  do
9:                           $v_{j,u,z}^k = \sum_{i=1}^N v_i^k p(z \mid x_i) p(x_i \mid u, x_j)$ 
10:                     endfor
11:                 endfor
12:             endfor
13:         endfor
14:         for all control actions  $u$  do
15:             for all  $k(1), \dots, k(M) = (1, \dots, 1)$  to  $(|\Upsilon|, \dots, |\Upsilon|)$  do
16:                 for  $i = 1$  to  $N$  do
17:                      $v'_i = \gamma \left[ r(x_i, u) + \sum_z v_{u,z,i}^{k(z)} \right]$ 
18:                 endfor
19:                 add  $(u; v'_1, \dots, v'_N)$  to  $\Upsilon'$ 
20:             endfor
21:         endfor
22:         optional: prune  $\Upsilon'$ 
23:          $\Upsilon = \Upsilon'$ 
24:     endfor
25:     return  $\Upsilon$ 

```

$$|A| = 3$$

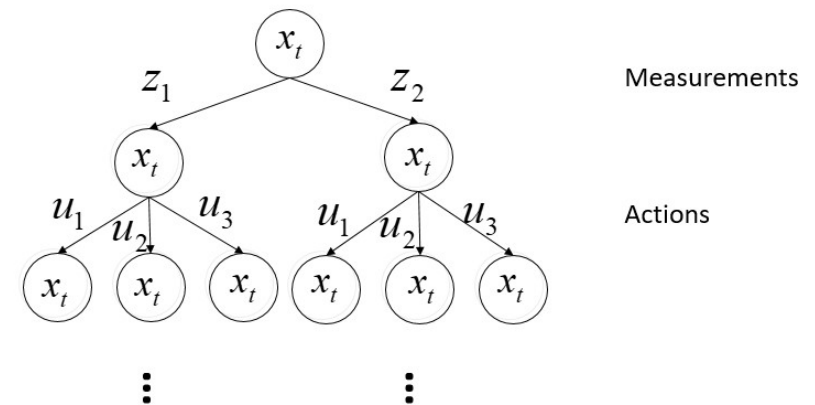
$$|Z| = 2$$

$$d = 11$$

$$3^{O(2^{\{d-1\}})}$$

$$= 3^{\{1024\}} \dots @ @ "$$

$ A $: actions
$ Z $: observations
d	: depth
$ A ^{O(Z ^{d-1})}$	paths



POMDP

- Why Pruning is Essential?
- Each **update introduces additional linear components** to V .
- Each **measurement squares the number of linear components**.
- Thus, an un-pruned value function for $T=20$ includes more than $10^{547,864}$ linear functions.
- At $T=30$ we have $10^{561,012,337}$ linear functions.
- The pruned value functions at $T=20$, in comparison, contains only 12 linear components.
- The combinatorial explosion of linear components in the value function are the major reason why **POMDPs are impractical for most applications**.

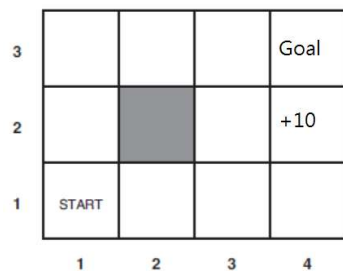
POMDP Summary

- POMDPs compute the optimal action in partially observable, stochastic domains.
- For finite horizon problems, the resulting value functions are piecewise linear and convex.
- In each iteration the number of linear constraints grows exponentially.
- POMDPs so far have only been applied successfully to very small state spaces with small numbers of possible observations and actions.

Conclusions

- LRTA*

Deterministic action



$$s, a \rightarrow s'$$

L2: Uninformed search

L3: Heuristic search (LRTA*)

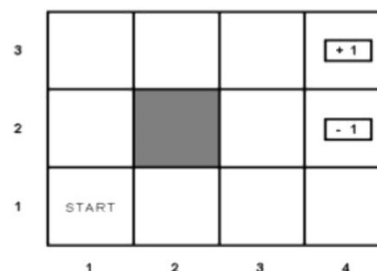
L4: Adversarial search

L5: Bayes theorem

L6: Bayes theorem over time

MDP (RL)

Probabilistic actions



$$P(s'|s, a)$$

L7: MDP

L9: Reinforcement learning

L10: GP and LWPR

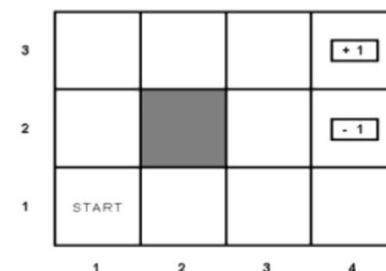
L11: Naïve Bayes and Perceptron

L12: Adaboost

L13: Deep learning and DRL

POMDP

Probabilistic actions and states

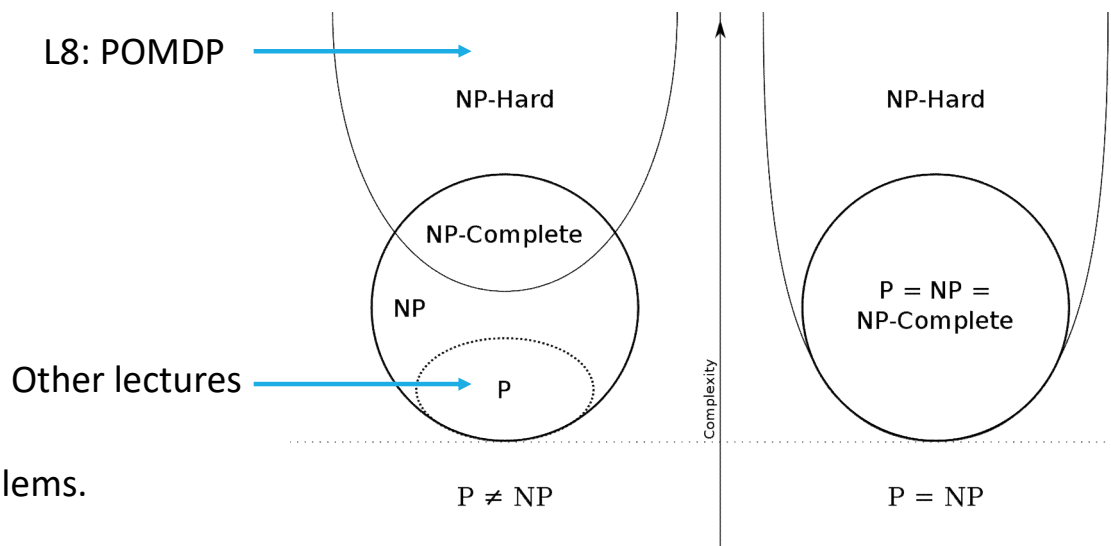


$$P(s'|s, a), P(s)$$

L8: POMDP

Appendix – NP-hard problems

- NP-hardness (non-deterministic polynomial-time hardness)
- NP-complete: Class of decision problems which contains the hardest problems in NP. Each NP-complete problem has to be in NP.



In MAI, we will try to solve P problems.

In L8, we will face POMDP, one of NP-hard problems.

<https://en.wikipedia.org/wiki/NP-hardness>

Q&A

