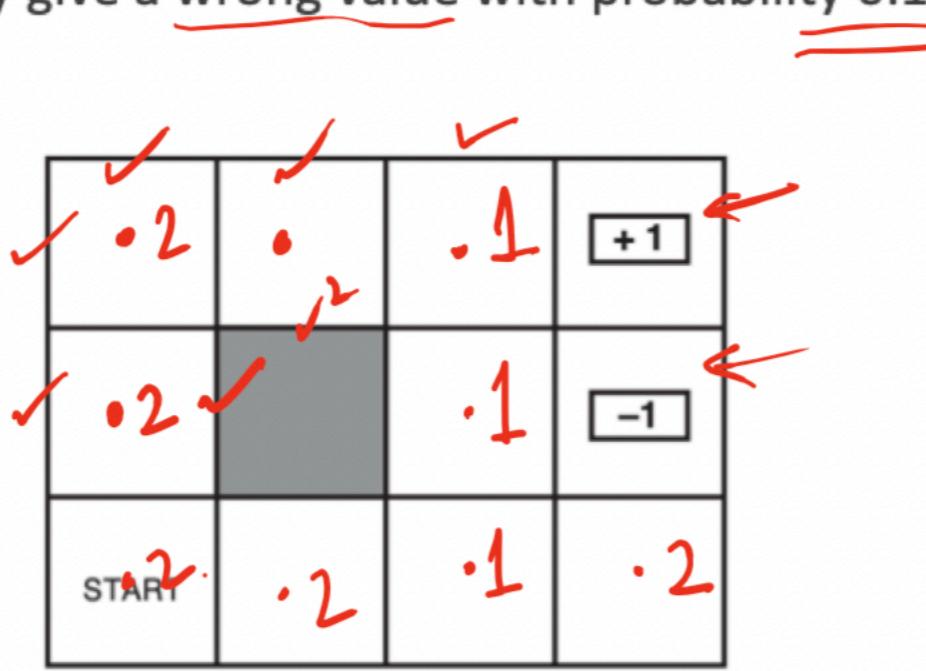


35

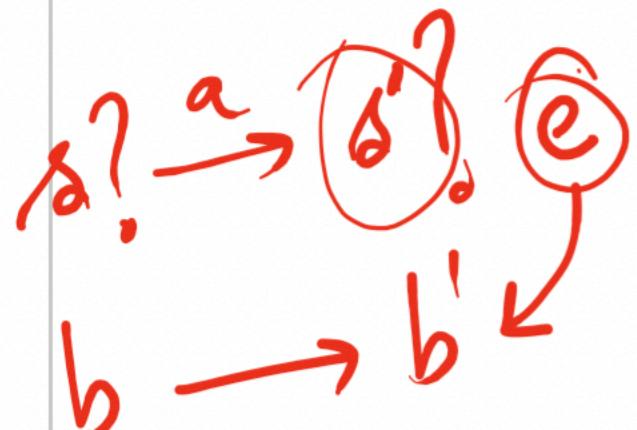
POMDP  $\xrightarrow{\text{convert}}$  MDP

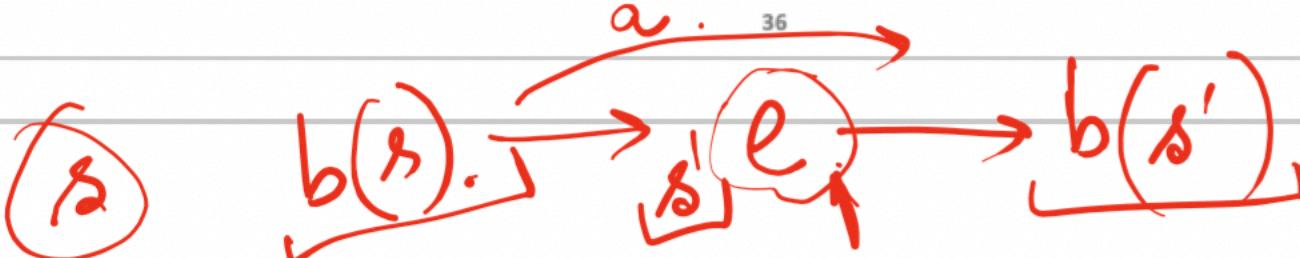
transition  
 $P(b'|b, a, e)$

- The  $4 \times 3$  world can make use of a noisy sensor that can measure the number of adjacent walls, but may give a wrong value with probability 0.1



$$\begin{array}{c}
 P(s'|b, a) \\
 \hline
 \begin{array}{c} 4 \times 3. \\ b(b) \\ P(e|s') \end{array} & P(b'|b, a, e)
 \end{array}$$





- The agent can compute its current belief state as a conditional probability distribution over the given sequence of percepts and actions so far.
- If  $b(s)$  was the previous belief state and the agent does action  $a$  and then perceives evidence  $e$  then the new belief state  $b'$  is given by

$$b'(s') = \alpha P(e | s') \sum_s P(s' | s, a) b(s)$$

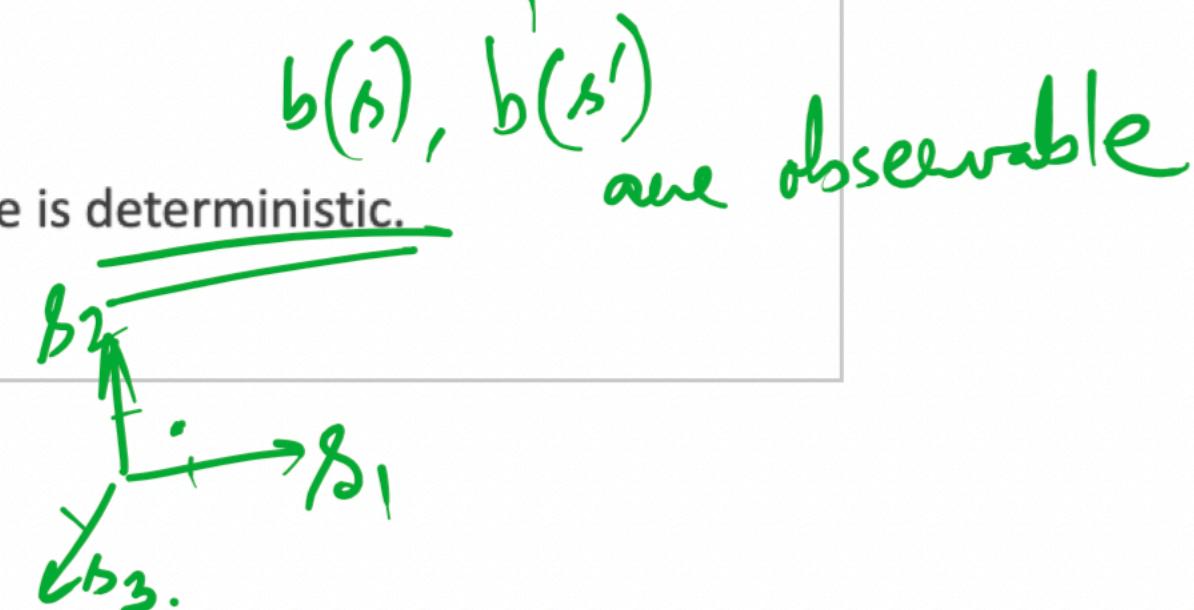
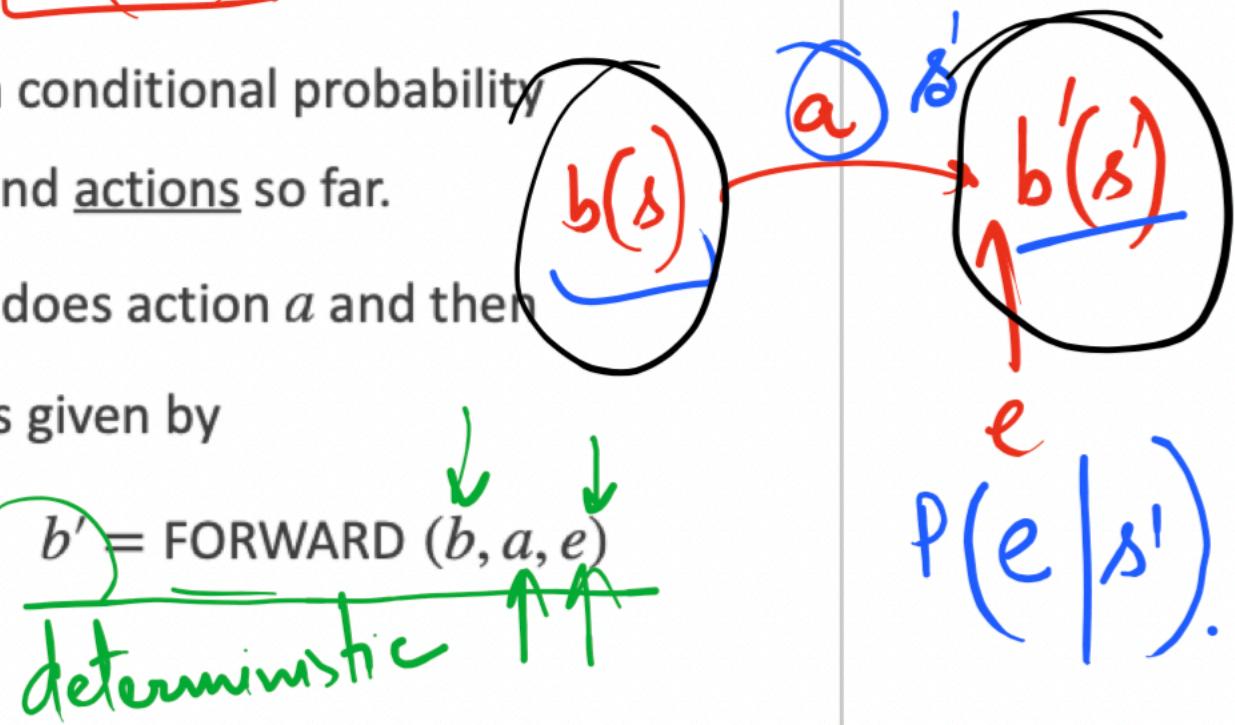
unnormalized

$$b'_{un}(s') = P(e | s') \sum_s P(s' | s, a) b(s)$$

- Given an evidence, the transition of the belief state is deterministic.

$b, b'$  are continuous.

37



- The optimal action depends only on the agent's current belief state and is given

by  $\pi^*(b)$ .

*Policy to be learned in the belief state*

- The decision cycle of a POMDP agent is as follows:

- Given the current belief state, execute the action  $a = \pi^*(b)$

- Receive percept  $e$

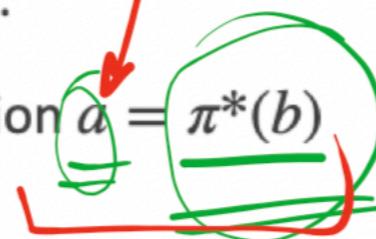
- Set the current belief state to  $\text{FORWARD}(b, a, e)$

- Repeat

$\pi^*(b) \rightarrow \text{plan}$  } *seq of action*

$b'$  is deterministic

*Value iteration*



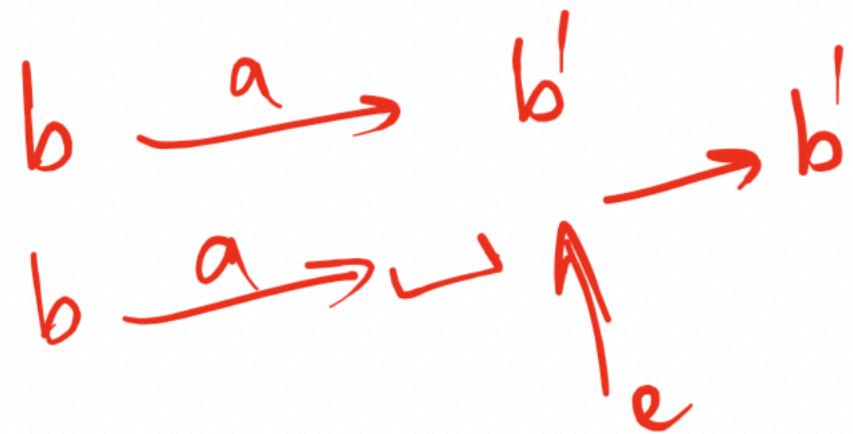
*first action of the plan*

*Compute ?*

*Utility to the physical states*

*belief states*

- The POMDP belief state space is continuous because the belief state is a probability distribution
- An action changes the belief state and not just the physical state.



Not the physical  $P(s' | s, a)$  POMDP

## Transition Model for the belief state space

- Without the evidence  $e$ , the transition of the belief state space is non-deterministic

$$\begin{aligned}
 P(b' | b, a) &= P(b' | a, b) = \sum_e P(b' | e, a, b) P(e | a, b) \\
 P(e | a, b) &= \sum_{s'} P(e | a, s', b) P(s' | a, b) \\
 &= \sum_{s'} P(e | s') P(s' | a, b) \\
 &= \sum_{s'} P(e | s') \sum_s P(s' | s, a) b(s) = \sum_{s'} b'_{un}(s')
 \end{aligned}$$

noisy sensor model.

## Transition Model for the belief state space

MDP

- The probability of reaching  $b'$  from  $b$ , given action  $a$  is  $P(b' | b, a)$

$$\begin{aligned} P(b' | b, a) &= \sum_e P(b' | e, a, b) P(e | a, b) \\ &= \sum_e P(b' | e, a, b) \sum_{s'} P(e | s') \sum_s P(s' | a) b(s) \end{aligned}$$

- Where  $P(b' | e, a, b) = 1$  if  $b' = \text{Forward}(b, a, e)$  and 0 otherwise.

$$= \sum_e P(b' | e, a, b) \sum_{s'} b'_{un}(s')$$

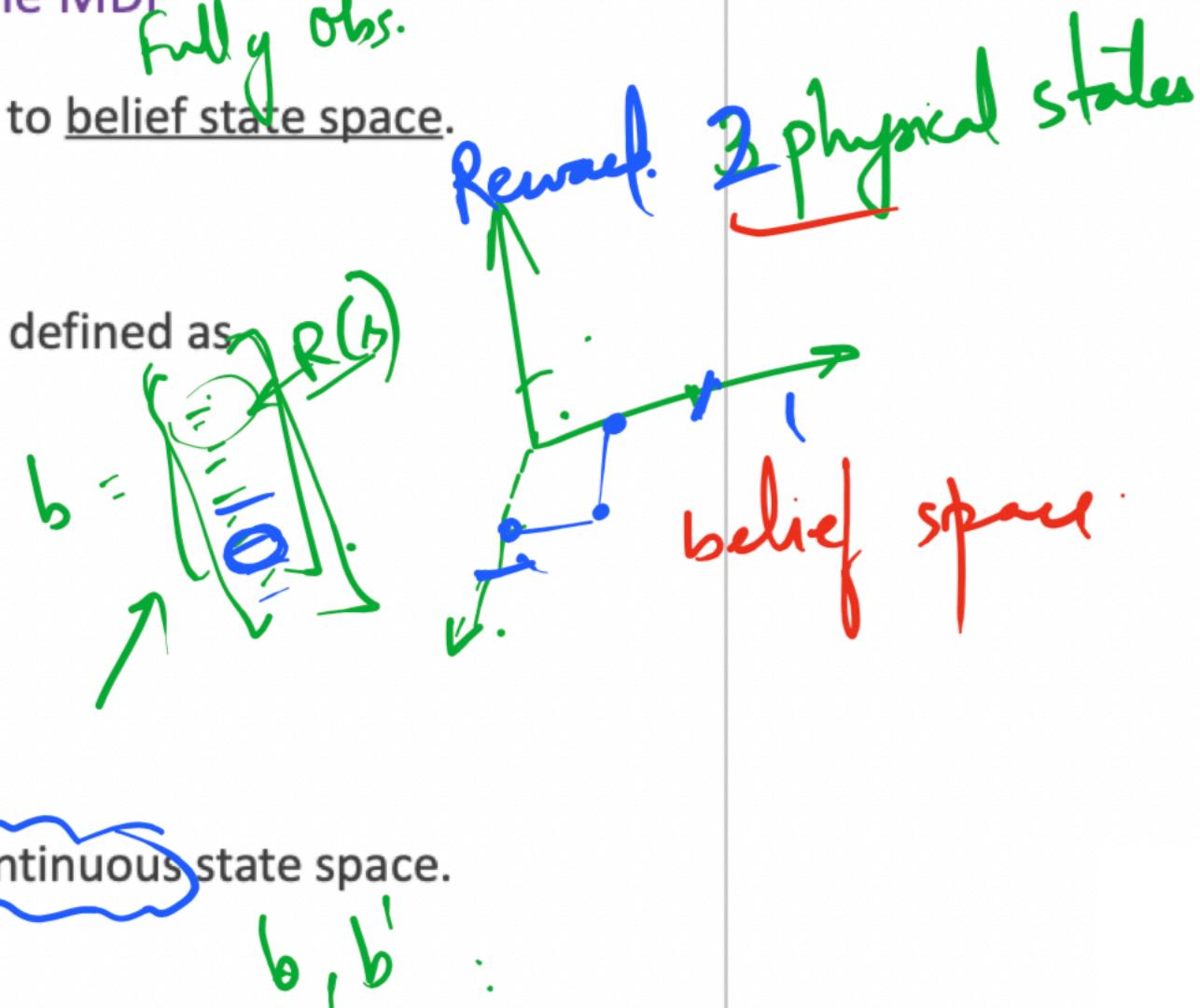
## Transforming partially observable MDP to observable MDP

Partially obs.      Fully obs.

- We transform the problem from physical state space to belief state space.
- Transition model in the belief state space  $P(b'|b, a)$
- The reward function on the belief state space can be defined as

$$\rho(b) = \sum_s b(s) R(s)$$

- The optimal policy for this MDP is  $\pi^*(b)$ .
- The belief state is always observable to the agent.
- BUT, this observable MDP has a high-dimensional continuous state space.



Value Iteration for POMDPs

$$\pi(b)$$

Best action

$$\pi^*(b)$$

Non trivial

Plan

plan

A1

Percept

$$p_1 \rightarrow A2'$$

$$p_2 \rightarrow A2''$$

$$p_3 \rightarrow A2'''$$

Percept

$$E'$$

$$E''$$

action

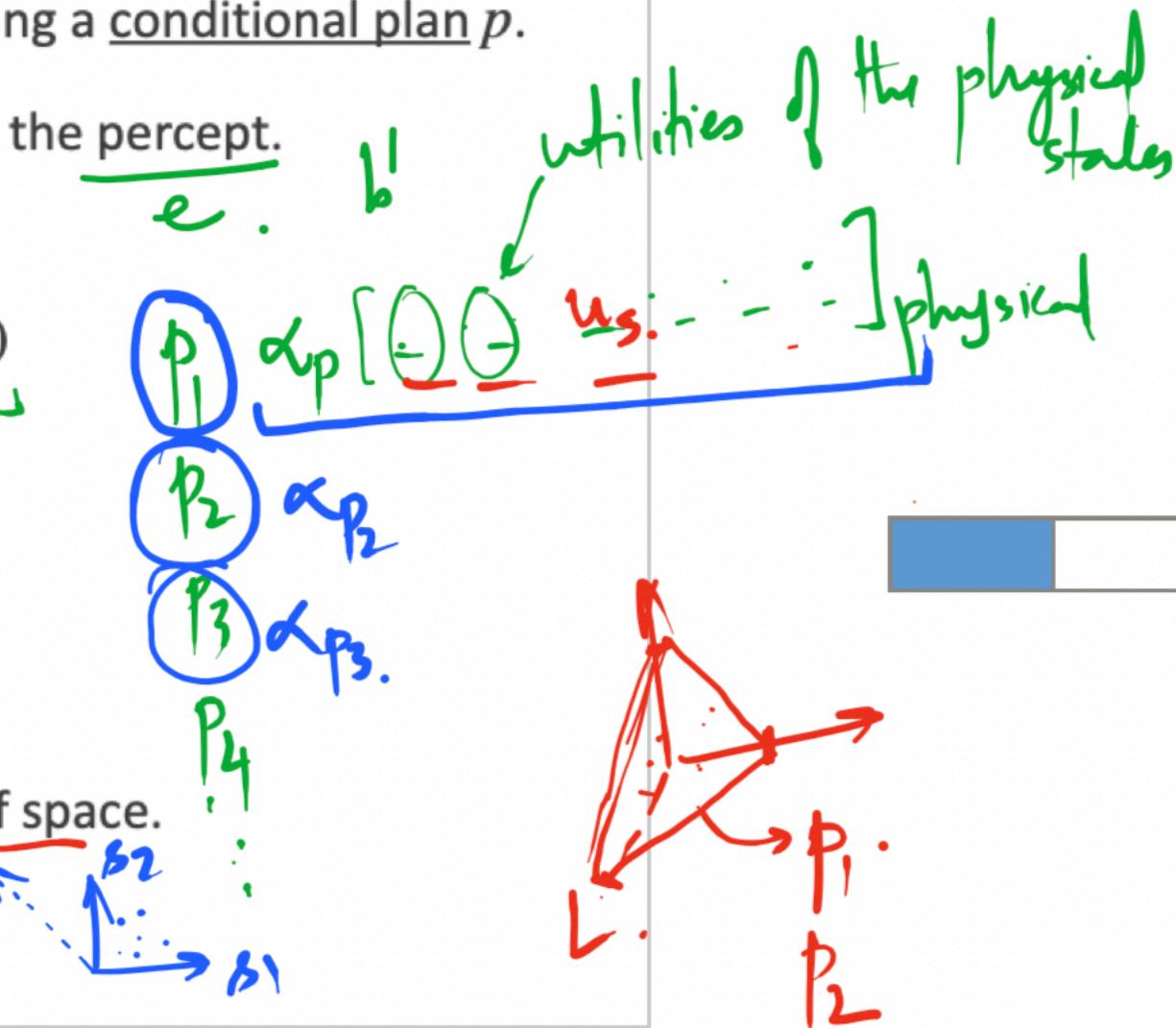
## What is $\pi^*(b)$ ?

action Plan

- The optimal policy for a belief state space is about choosing a conditional plan  $p$ .
- The belief state is updated based on the action  $a$  as well as the percept  $e$ .
- Let the utility of executing a fixed conditional plan  $p$  starting in physical state  $s$  is  $\alpha_p(s)$
- The expected utility of executing  $p$  in belief state  $b$  is

$$\text{expected utility} \rightarrow \sum_s b(s) \alpha_p(s) \quad \text{or} \quad b \cdot \alpha_p$$

The expected utility of plan  $p$  is a hyperplane in the belief space.

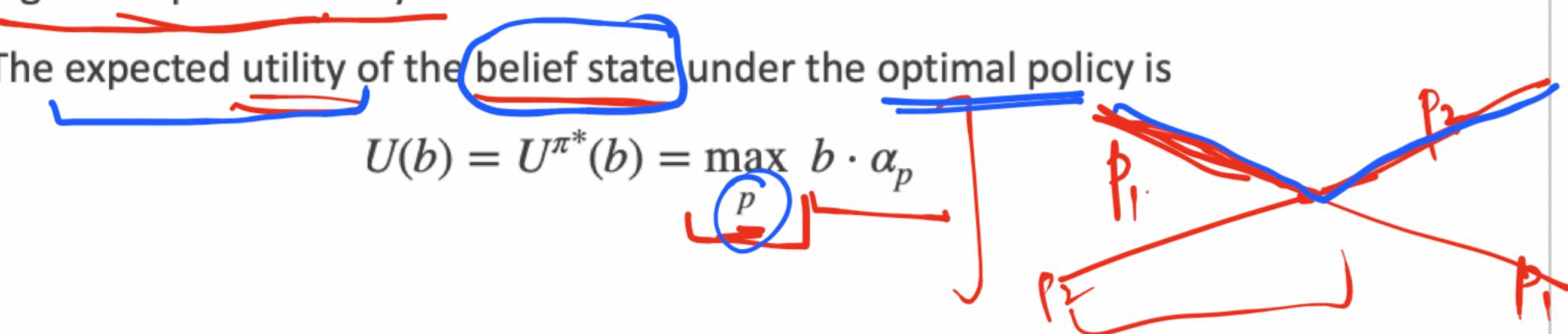


## Optimal Policy $\pi^*(b)$

- The optimal policy in the belief state space chooses a conditional plan with highest expected utility.
- The expected utility of the belief state under the optimal policy is

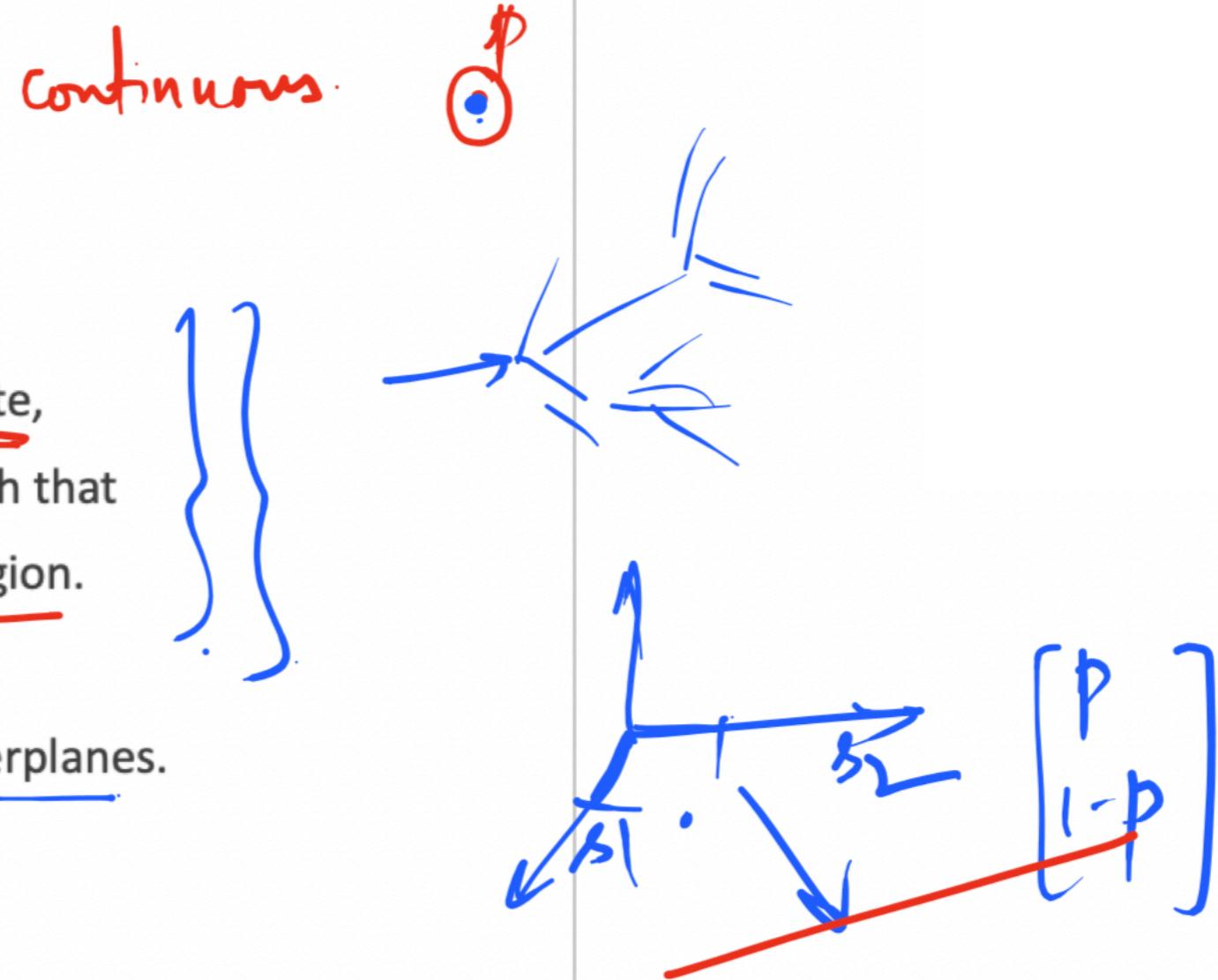
$$U(b) = U^{\pi^*}(b) = \max_p b \cdot \alpha_p$$

$b$  exp.  
 $U(b)$  best plan.



## Implication of a continuous belief state in POMDPs

- We can expect the agent to execute plan  $p$  in belief states that are very close to  $b$ .
- If the number of conditional plans being considered is finite, then the belief state space can be divided into regions such that each region has a particular plan that is optimal in that region.
- The utility function is the maximum of a collection of hyperplanes.

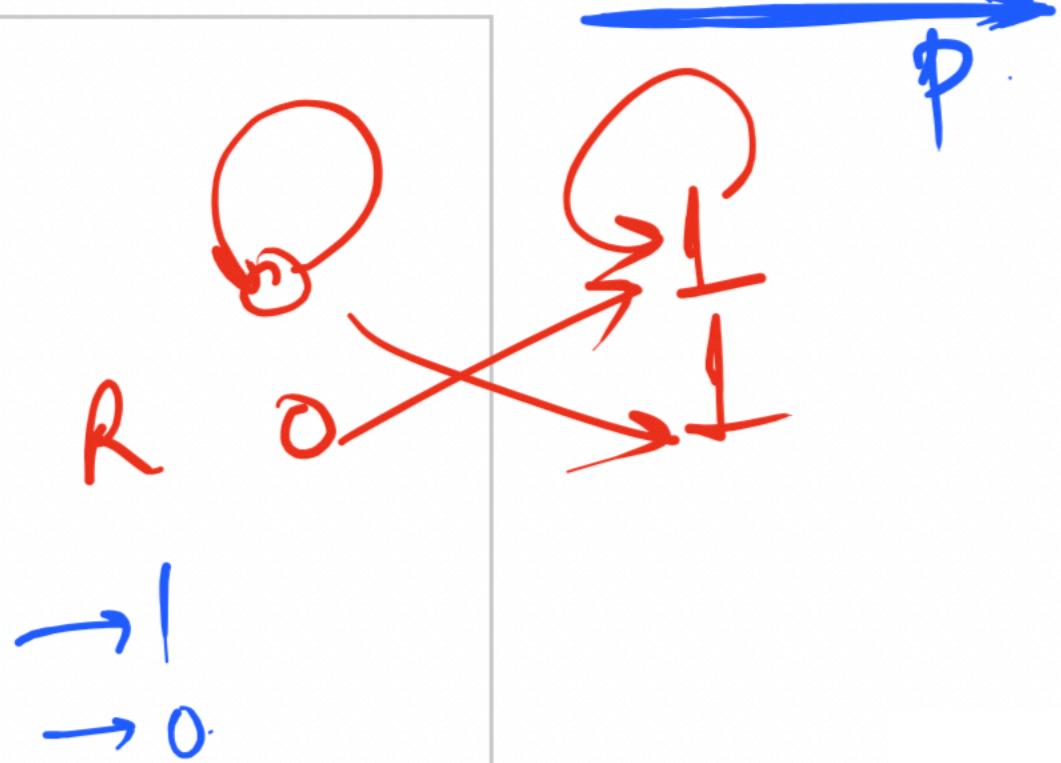


## Example: A two-state world

- State with label 0 carries reward  $R(0) = 0$
- State with label 1 carries reward  $R(1) = 1$
- Action 1: Stay  
Stay put with probability  $0.9$
- Action 2: Go  
Go switches to the other state with probability  $0.9$
- Sensor reports the correct ~~the~~ state with probability 0.6



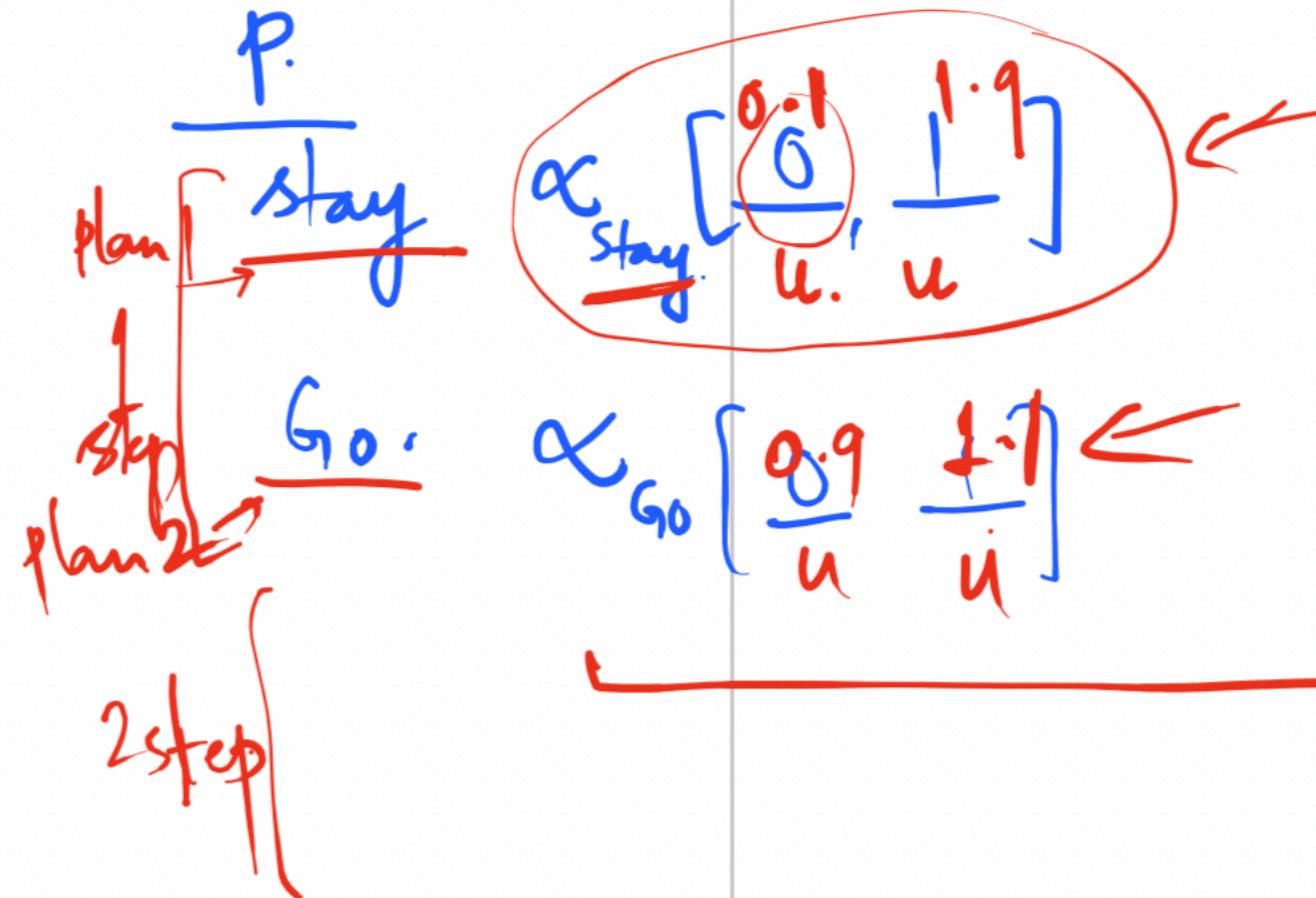
$$\begin{array}{ccc} 0 & \xrightarrow{\hspace{1cm}} & 1 \\ 1 & \xrightarrow{\hspace{1cm}} & 0 \end{array}$$



- Next we should assign expected utilities  $\alpha_p$  to the states {0,1} based on 1-step plans

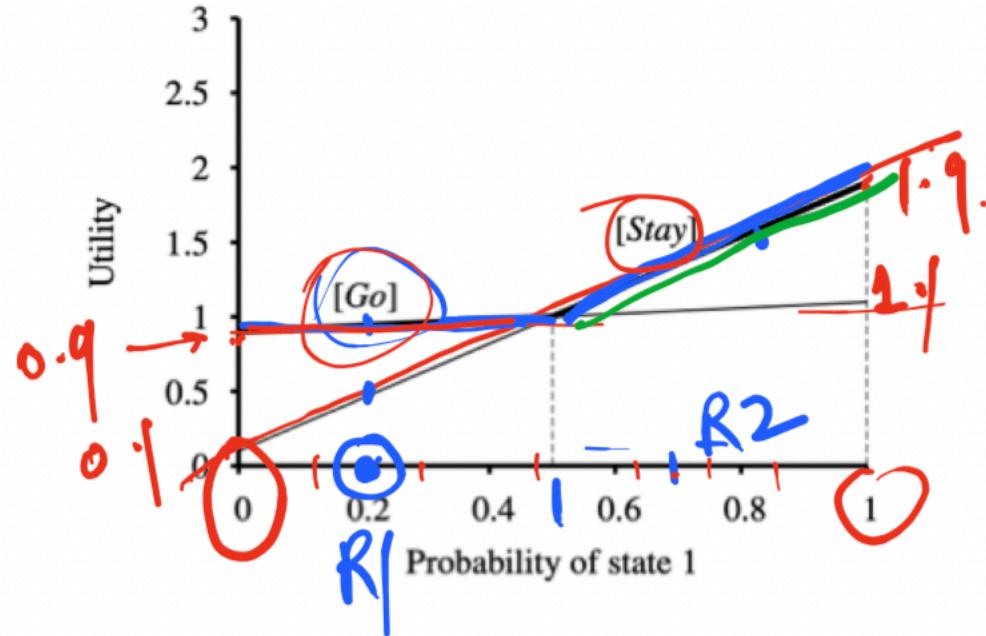
$$\begin{aligned}\alpha_{[Stay]}(0) &= R(0) + \gamma(0.9R(0) + 0.1R(1)) = 0.1 \\ \alpha_{[Stay]}(1) &= R(1) + \gamma(0.9R(1) + 0.1R(0)) = 1.9 \\ \alpha_{[Go]}(0) &= R(0) + \gamma(0.9R(1) + 0.1R(0)) = 0.9 \\ \alpha_{[Go]}(1) &= R(1) + \gamma(0.9R(0) + 0.1R(1)) = 1.1\end{aligned}$$

$R(s) + \gamma \sum R(s')$



- The hyperplanes for  $b \cdot \alpha_{[Stay]}$  and  $b \cdot \alpha_{[Go]}$ .

$$\alpha_{[Stay]} \begin{bmatrix} 0.1 \\ 1.9 \end{bmatrix}$$

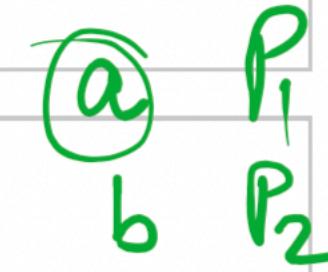
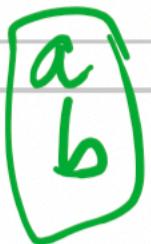


- The bold line represents the utility function for the finite horizon that allows only one action.

2 states

$$\begin{bmatrix} - \\ - \end{bmatrix} \xrightarrow{\quad} \begin{bmatrix} 0.1 \\ 1.9 \end{bmatrix}$$

+ 1.0

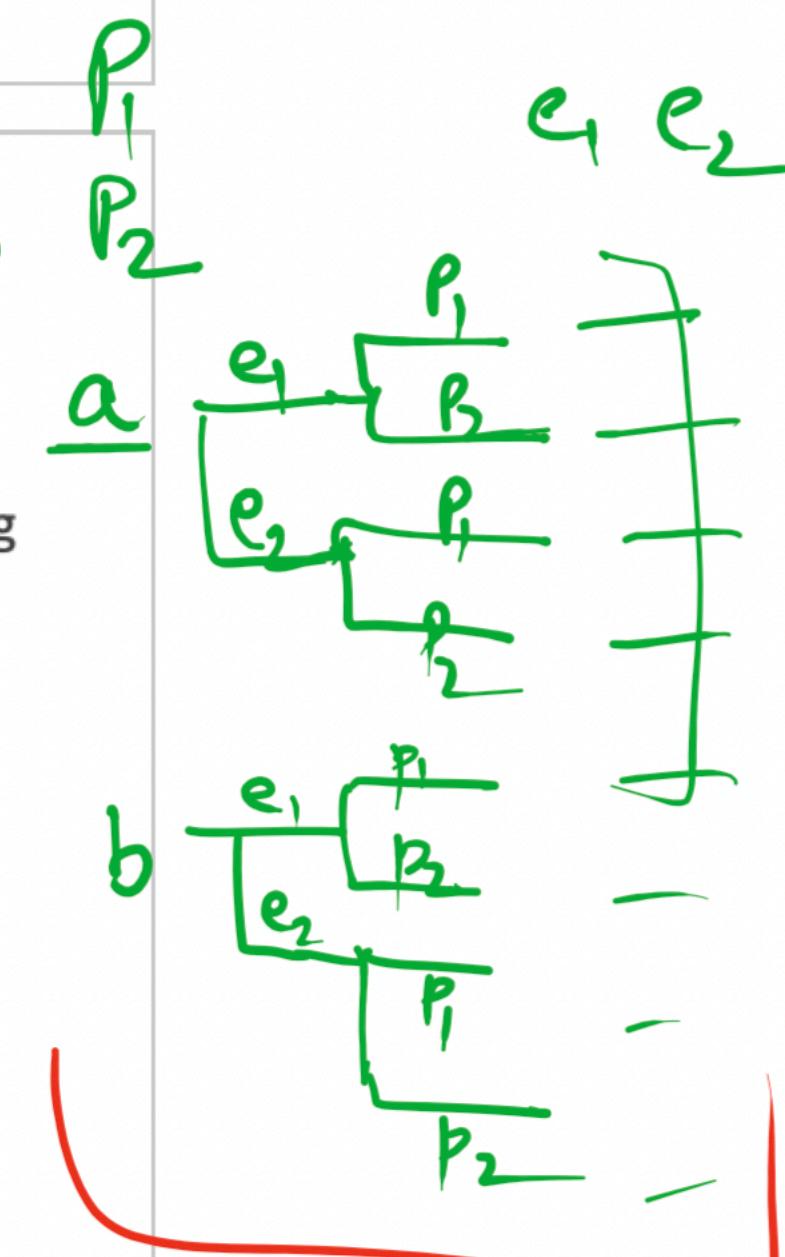


- We extend the horizon one step at a time and compute the utility for the physical state  $s$
- That is, we now compute utilities for conditional plans of depth 2 by considering
  - each possible first action,
  - each possible subsequent percept, and then
  - each way of choosing a depth-1 plan to execute for each percept

Utilities?

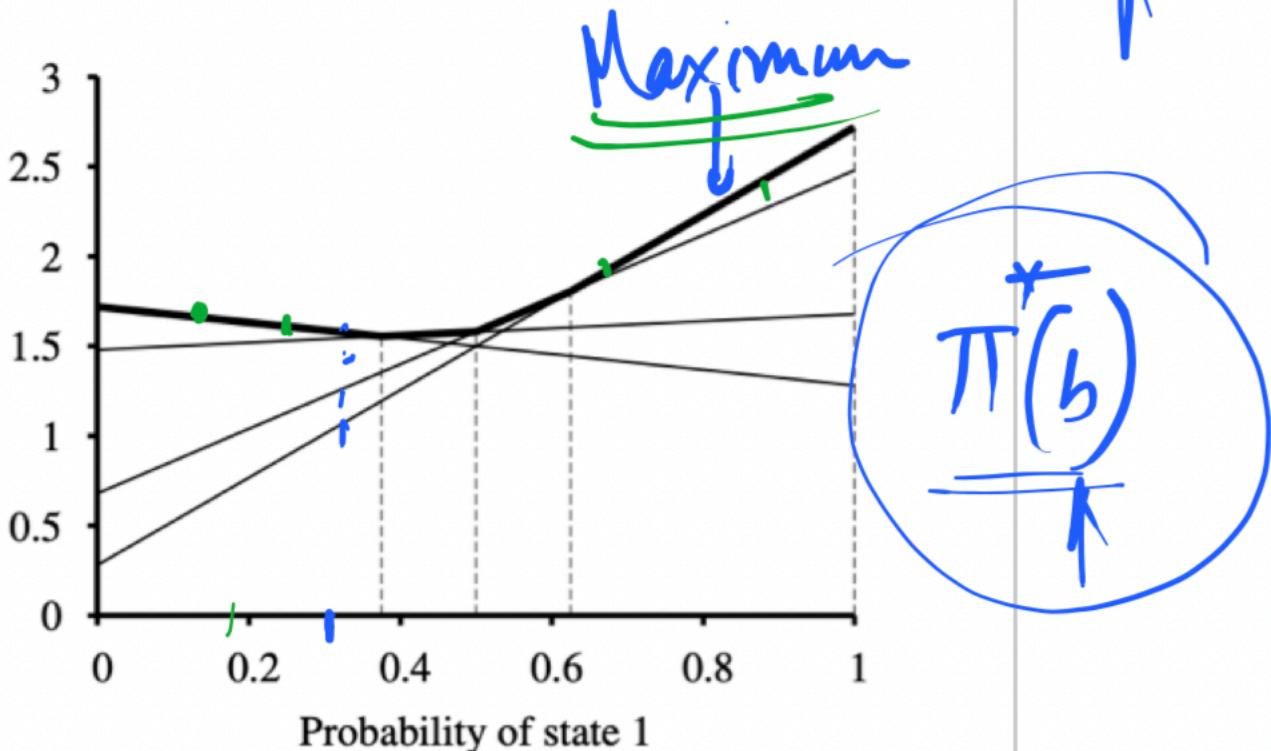
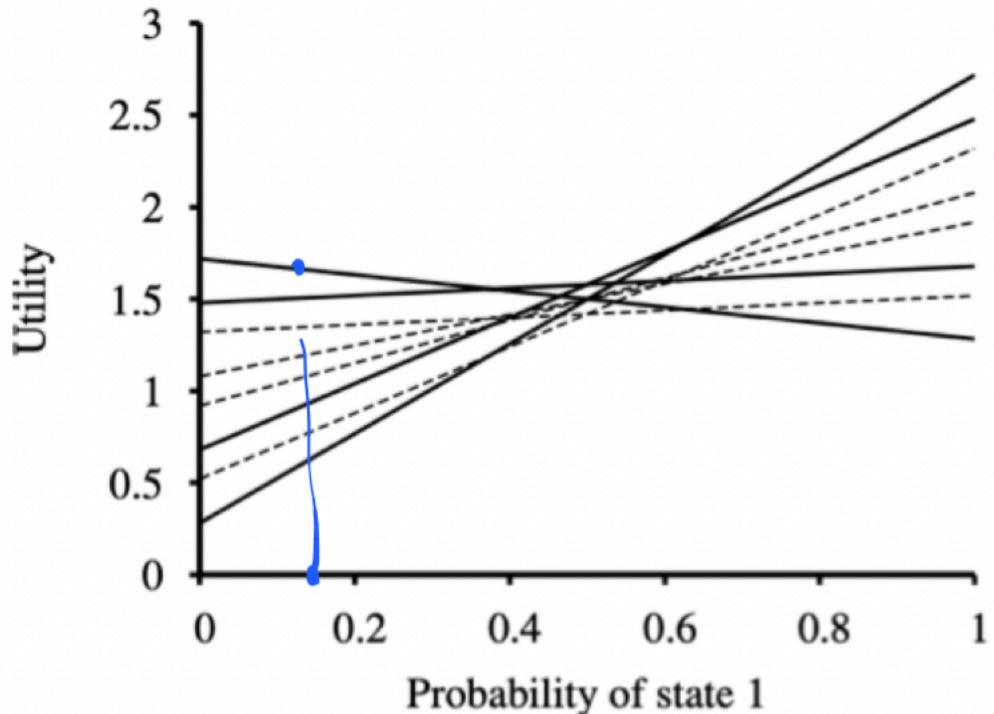
[Stay; **if** Percept = 0 **then** Stay **else** Stay]  
 [Stay; **if** Percept = 0 **then** Stay **else** Go] ...

$\alpha_{P_{new}}$



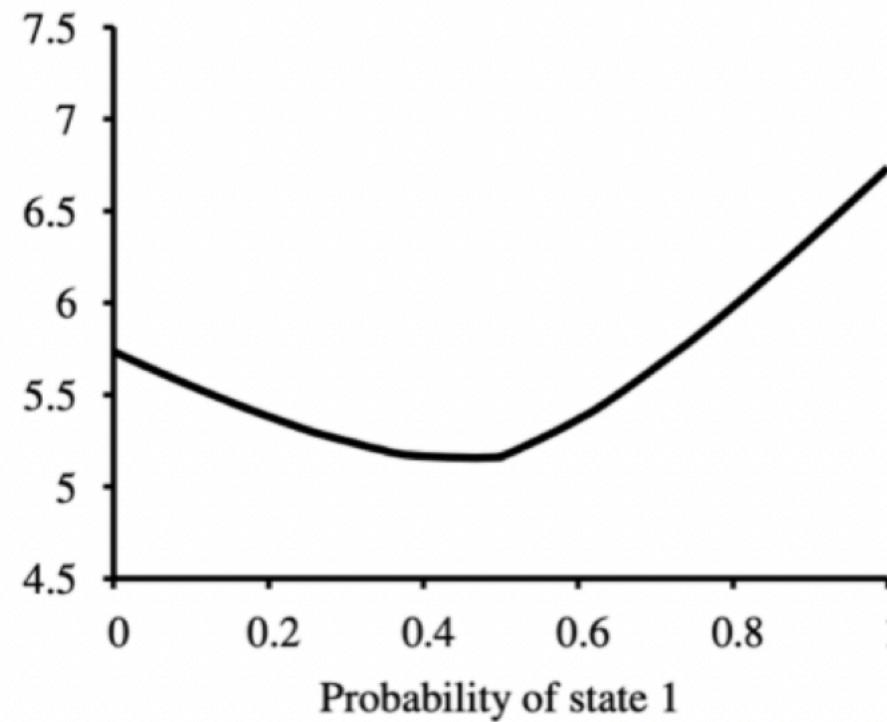
2 steps.

- There are 8 distinct depth 2 plans and their utility for physical state  $s$



- The dashed plans are suboptimal across the entire belief space

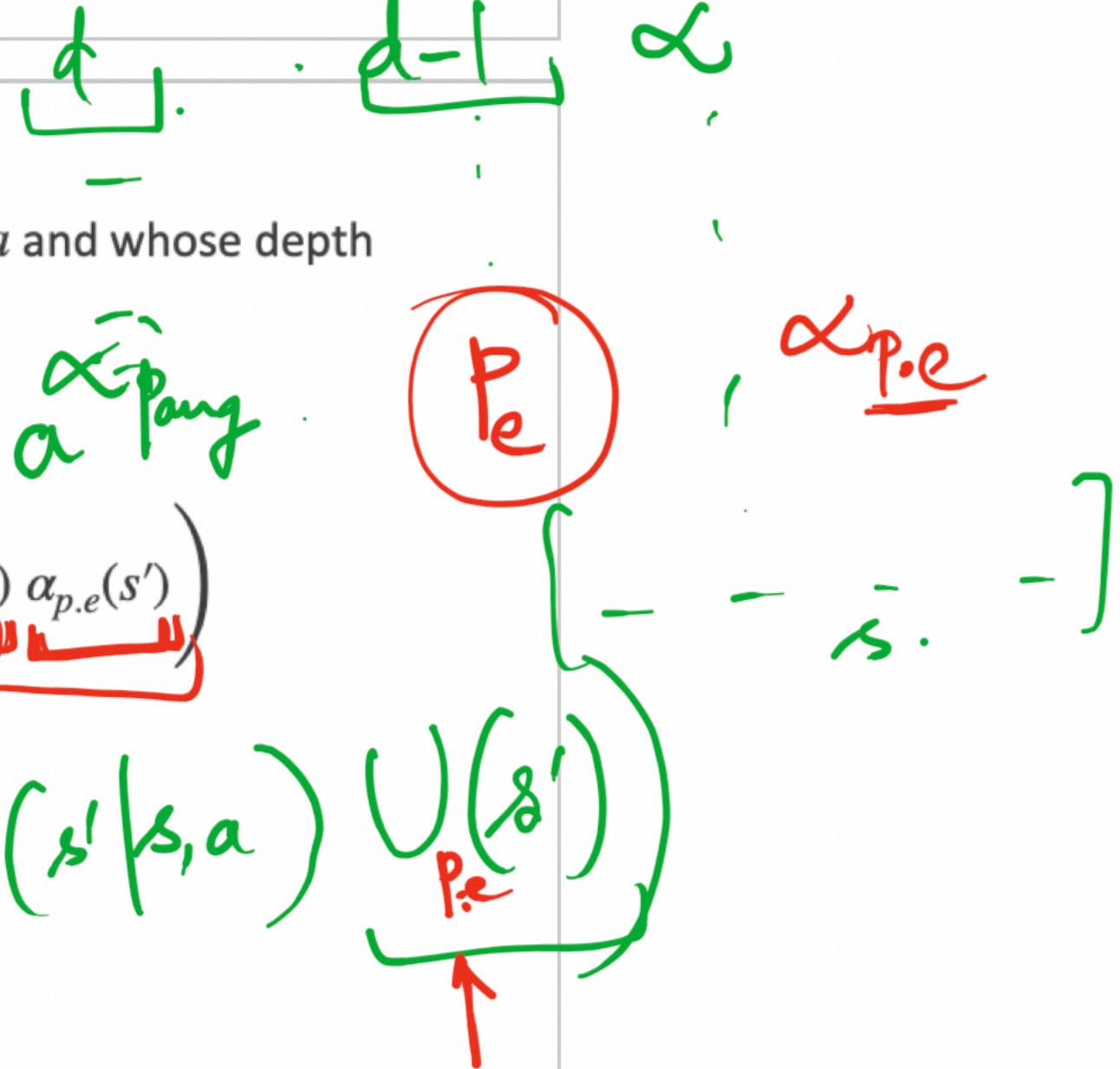
- If we extend the horizon to to 8-step plan the utility function is



Piece-wise constant,  
Convex

## Backing up the utility values for a depth $d$ plan

- Let  $p$  be a depth  $d$  conditional plan whose initial action is  $a$  and whose depth  $d - 1$  sub plan for percept  $e$  is  $p \cdot e$



$$\alpha_p(s) = R(s) + \gamma \left( \sum_{s'} P(s'|s,a) \sum_e P(e|s') \alpha_{p.e}(s') \right)$$

$$\alpha_p(s) = R(s) + \gamma \left( \sum_{s'} P(s'|s,a) U(s') \right)$$

## Value Iteration Algorithm

**function** POMDP-VALUE-ITERATION(*pomdp*,  $\epsilon$ ) **returns** a utility function

**inputs:** *pomdp*, a POMDP with states  $S$ , actions  $A(s)$ , transition model  $P(s' | s, a)$ ,  
sensor model  $P(e | s)$ , rewards  $R(s)$ , discount  $\gamma$   
 $\epsilon$ , the maximum error allowed in the utility of any state

**local variables:**  $U$ ,  $U'$ , sets of plans  $p$  with associated utility vectors  $\alpha_p$

$U' \leftarrow$  a set containing just the empty plan  $[]$ , with  $\alpha_{[]} (s) = R(s)$

**repeat**

$U \leftarrow U'$

$U' \leftarrow$  the set of all plans consisting of an action and, for each possible next percept,  
a plan in  $U$  with utility vectors computed according to

$U' \leftarrow$  REMOVE-DOMINATED-PLANS( $U'$ )

$$\alpha_p(s) = R(s) + \gamma \left( \sum_{s'} P(s' | s, a) \sum_e P(e | s') \alpha_{p,e}(s') \right)$$

**until** MAX-DIFFERENCE( $U$ ,  $U'$ )  $< \epsilon(1 - \gamma)/\gamma$

**return**  $U$

