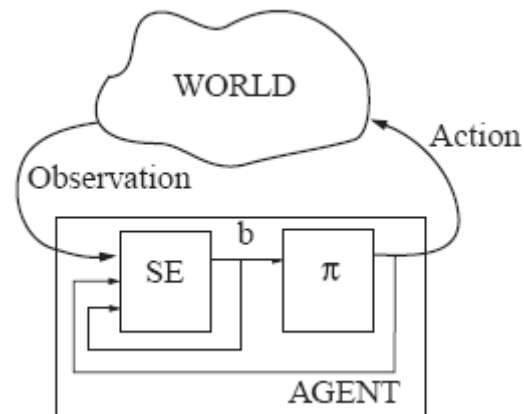# POMDPs

- MDPs policy: to find a mapping from states to actions

- POMDPs policy: to find a mapping from probability distributions (over states) to actions.

    - belief state: a probability distribution over states

    - belief space: the entire probability space, infinite

# POMDPs

■ **Partially Observable MDPs**

A partially observable Markov decision process can be described as a tuple $\langle \mathcal{S}, \mathcal{A}, T, R, \Omega, O \rangle$, where

- $\mathcal{S}$, $\mathcal{A}$, $T$, and $R$ describe a Markov decision process;
- $\Omega$ is a finite set of observations the agent can experience of its world; and
- $O : \mathcal{S} \times \mathcal{A} \rightarrow \Pi(\Omega)$ is the *observation function*, which gives, for each action and resulting state, a probability distribution over possible observations (we write $O(s', a, o)$ for the probability of making observation $o$ given that the agent took action $a$ and landed in state $s'$).

# POMDPs

- In POMDPs we apply the very same idea as in MDPs.

- **Since the state is not observable**, the agent has to **make its decisions based on** the belief state which is a **posterior distribution over states**.

- Let $b$ be the belief of the agent about the state under consideration.

- POMDPs compute a **value function over belief space**:

$$V_T(b) = \gamma \max_u \left[ r(b,u) + \int V_{T-1}(b')p(b' \mid u, b) \, db' \right]$$

# Problems

- Each belief is a probability distribution, thus, each value in a **POMDP is a function of an entire probability distribution**.

- **This is problematic, since probability distributions are continuous**.

- Additionally, we have to deal with the **huge complexity of belief spaces**.

- For **finite worlds** with finite state, action, and measurement spaces and finite horizons, however, we can **effectively represent the value functions by piecewise linear functions**.
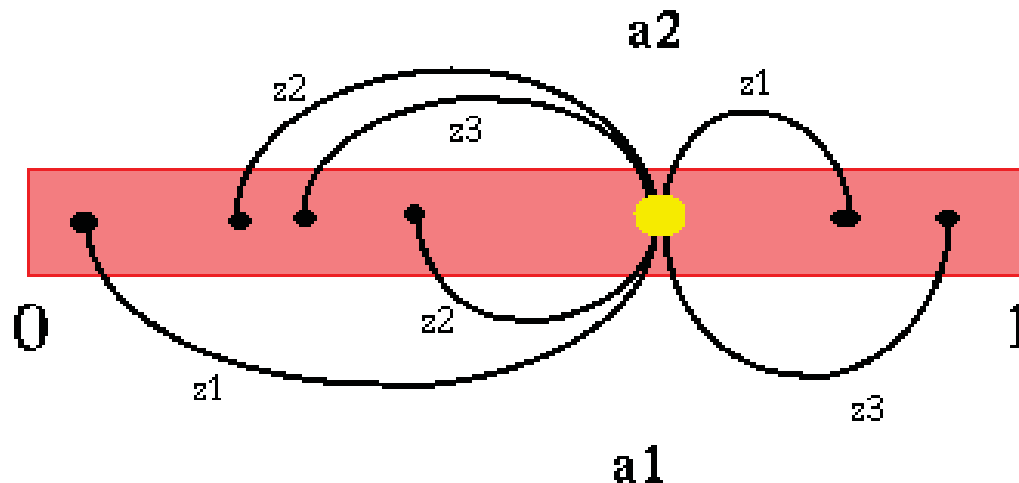
# A two state POMDP

- represent the belief state with a single number $p$.
- the entire space of belief states can be represented as a line segment.
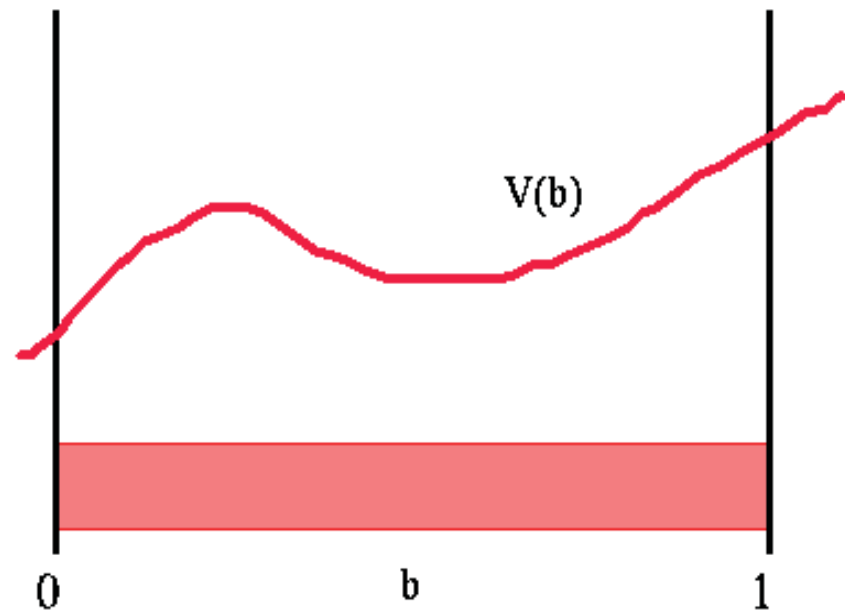
belief space for a 2 state POMDP

# belief state updating

- finite number of possible next belief states, given a belief state
  - a finite number of actions
  - a finite number of observations
- $b' = T(b' \mid b, a, z)$. Given $a$ and $z$, $b'$ is fully determined.

- the process of maintaining the belief state is Markovian: the next belief state depends only on the current belief state (and the current action and observation)

- we are now back to solving a MDP policy problem with some adaptations
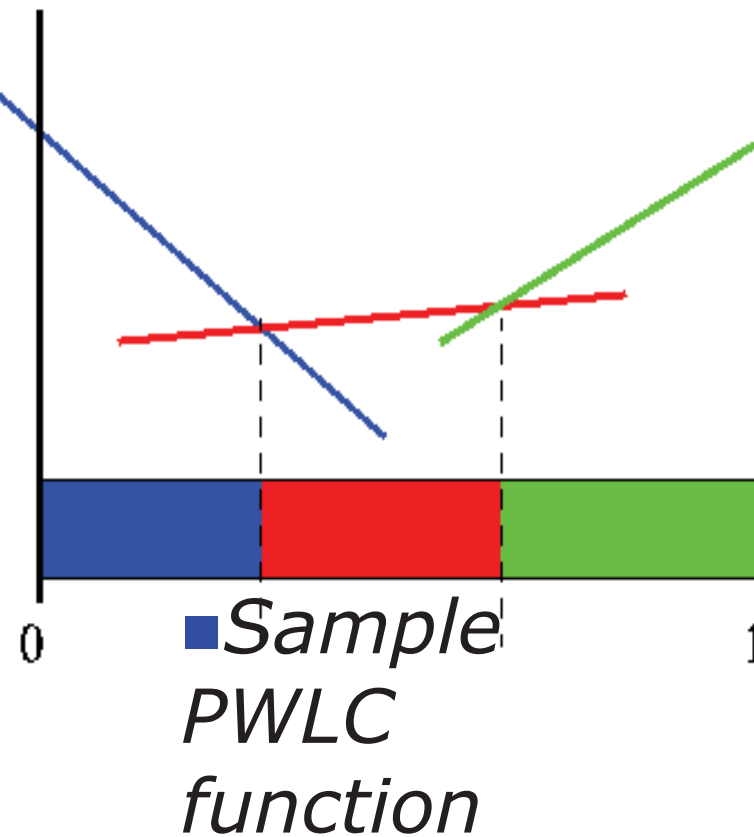
- **continuous space: value function is some arbitrary function**
  - b: belief space
  - V(b): value function

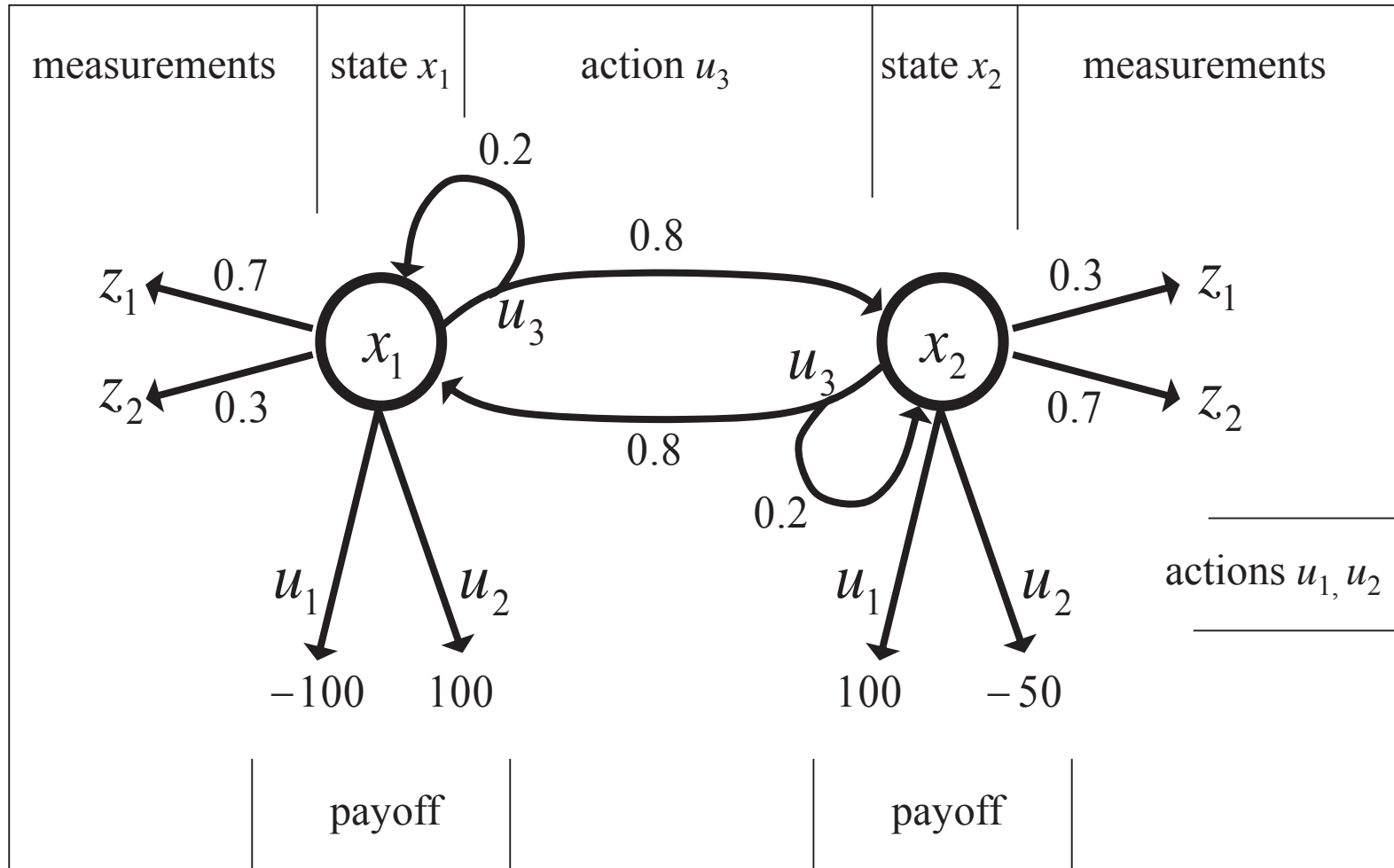- **Problem: how we can easily represent this value function?**



- *Value function over belief space*

57

Fortunately, the finite horizon value function is piecewise linear and convex (PWLC) for every horizon length.



0    ■*Sample PWLC function*    1

# An Illustrative Example

# The Parameters of the Example

- 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 $\gamma = 1$.

$$
\begin{aligned}
r(x_1, u_1) &= -100 & r(x_2, u_1) &= +100 \\
r(x_1, u_2) &= +100 & r(x_2, u_2) &= -50 \\
r(x_1, u_3) &= -1 & r(x_2, u_3) &= -1
\end{aligned}
$$

$$
\begin{aligned}
p(x_1'|x_1, u_3) &= 0.2 & p(x_2'|x_1, u_3) &= 0.8 \\
p(x_1'|x_2, u_3) &= 0.8 & p(z_2'|x_2, u_3) &= 0.2
\end{aligned}
$$

$$
\begin{aligned}
p(z_1|x_1) &= 0.7 & p(z_2|x_1) &= 0.3 \\
p(z_1|x_2) &= 0.3 & p(z_2|x_2) &= 0.7
\end{aligned}
$$

# Payoff in POMDPs

- 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 over all states**:

$$
\begin{aligned}
r(b, u) &= E_x[r(x, u)] \\
&= \int r(x, u) p(x) \, dx \\
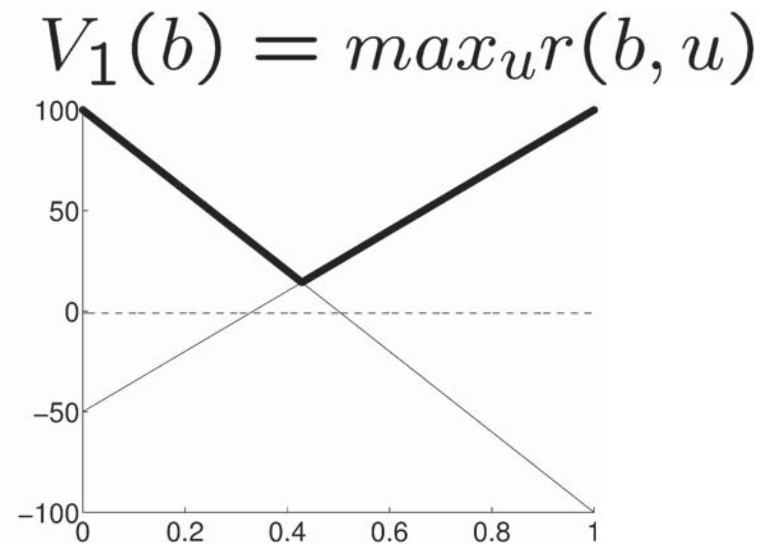&= p_1 \, r(x_1, u) + p_2 \, r(x_2, u)
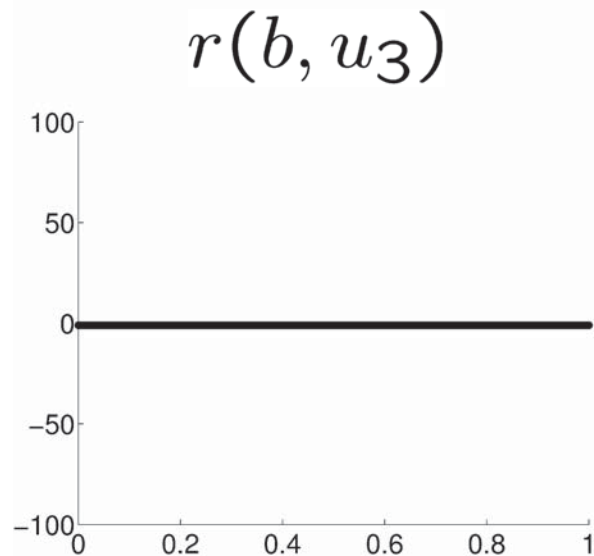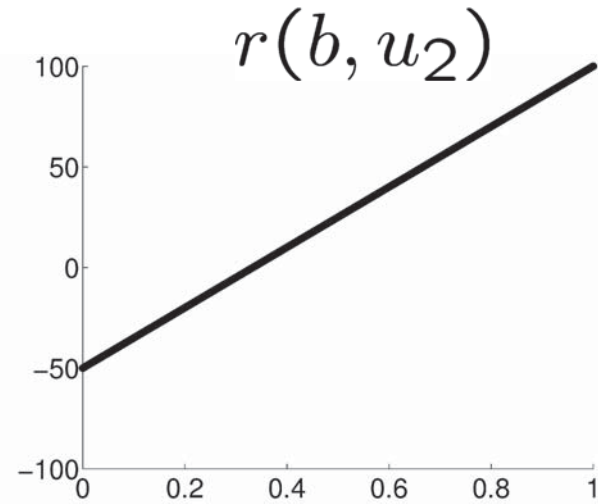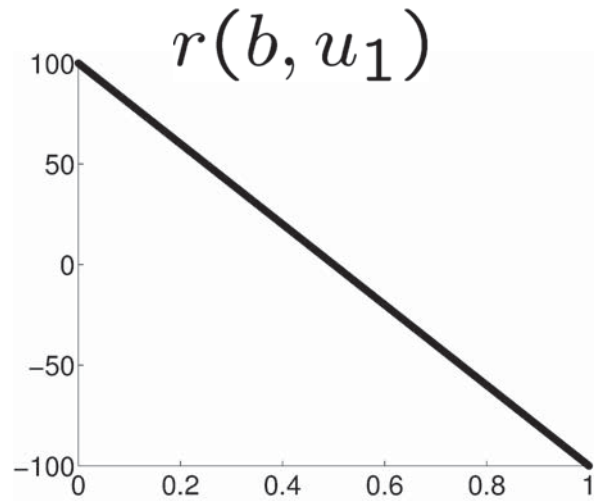\end{aligned}
$$

# Payoffs in Our Example (1)

- 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

$$r(b, u_1) = -100\, p_1 + 100\, p_2$$
$$= -100\, p_1 + 100\, (1 - p_1)$$

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

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

# Payoffs in Our Example (2)

# The Resulting Policy for T=1

- Given we have a finite POMDP with T=1, we would use $V_1(b)$ to determine the optimal policy.

- In our example, the optimal policy for T=1 is

$$\pi_1(b) \;=\; \begin{cases} u_1 & \text{if } p_1 \leq \frac{3}{7} \\[2ex] u_2 & \text{if } p_1 > \frac{3}{7} \end{cases}$$

- This is the upper thick graph in the diagram.

# Piecewise Linearity, Convexity

■ The resulting value function $V_1(b)$ is the maximum of the three functions at each point

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

$$= \max \left\{ \begin{array}{cc} -100\, p_1 & +100\, (1 - p_1) \\ 100\, p_1 & -50\, (1 - p_1) \\ & -1 \end{array} \right\}$$
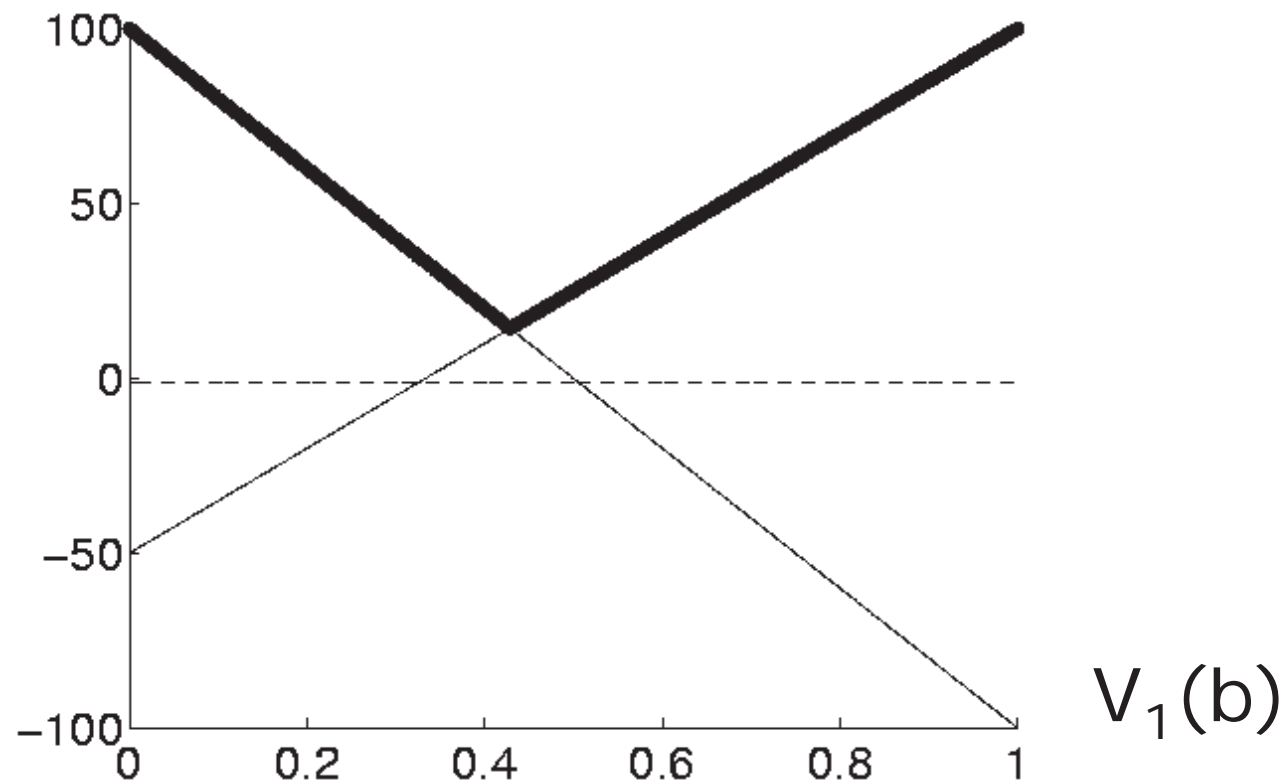
■ It is piecewise linear and convex.

# Pruning

- If we carefully consider $V_1(b)$, we see that only the first two components contribute.

- The third component can therefore safely be pruned away from $V_1(b)$.

$$V_1(b) = \max\left\{ \begin{array}{ll} -100\ p_1 & +100\ (1-p_1) \\ 100\ p_1 & -50\ (1-p_1) \end{array} \right\}$$

# Increasing the Time Horizon

- Assume the robot can make an observation before deciding on an action.
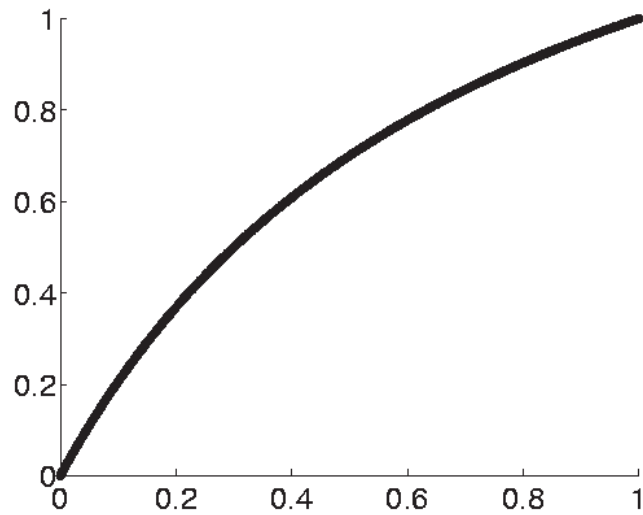


$V_1(b)$

# Increasing the Time Horizon

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

$$p'_1 = \frac{0.7\,p_1}{p(z_1)}$$
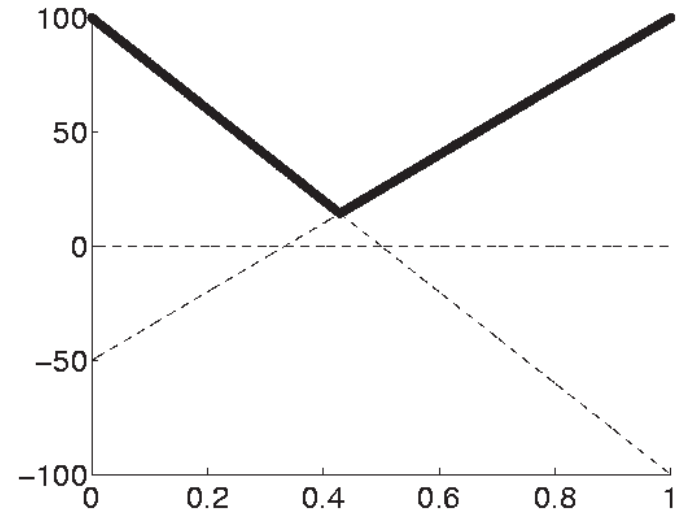
$$p'_2 = \frac{0.3(1 - p_1)}{p(z_1)}$$

$$p(z_1) = 0.7\,p_1 + 0.3(1 - p_1) = 0.4\,p_1 + 0.3$$
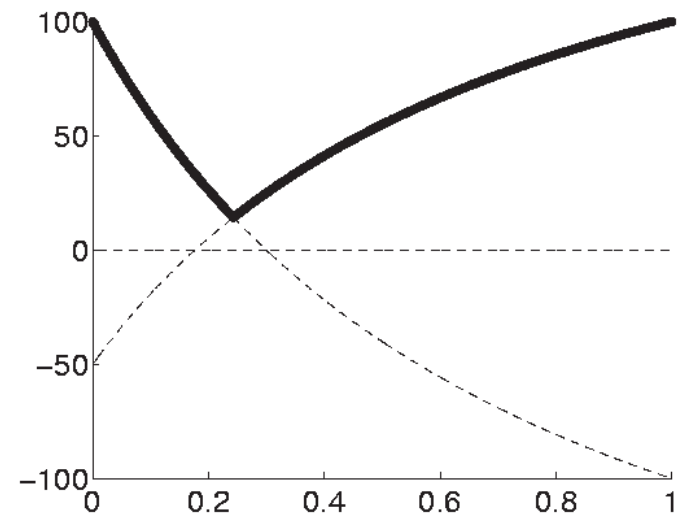
# Value Function

$V_1(b)$

$b'(b|z_1)$

$V_1(b|z_1)$

# Increasing the Time Horizon

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

$$V_1(b \mid z_1) \;=\; \max \left\{ \begin{array}{cc} -100 \cdot \frac{0.7\, p_1}{p(z_1)} & +100 \cdot \frac{0.3\,(1-p_1)}{p(z_1)} \\[2ex] 100 \cdot \frac{0.7\, p_1}{p(z_1)} & -50 \cdot \frac{0.3\,(1-p_1)}{p(z_1)} \end{array} \right\}$$

$$= \; \frac{1}{p(z_1)} \max \left\{ \begin{array}{cc} -70\, p_1 & +30\,(1-p_1) \\ 70\, p_1 & -15\,(1-p_1) \end{array} \right\}$$

# Expected Value after Measuring

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

$$\overline{V}_1(b) = E_z[V_1(b \mid z)] = \sum_{i=1}^{2} p(z_i) V_1(b \mid z_i)$$

$$= \sum_{i=1}^{2} p(z_i) V_1\left( \frac{p(z_i \mid x_1) p_1}{p(z_i)} \right)$$

$$= \sum_{i=1}^{2} V_1\left( p(z_i \mid x_1) p_1 \right)$$

# Expected Value after Measuring

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

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

# Resulting Value Function

- The four possible combinations yield the following function which then can be simplified and pruned.
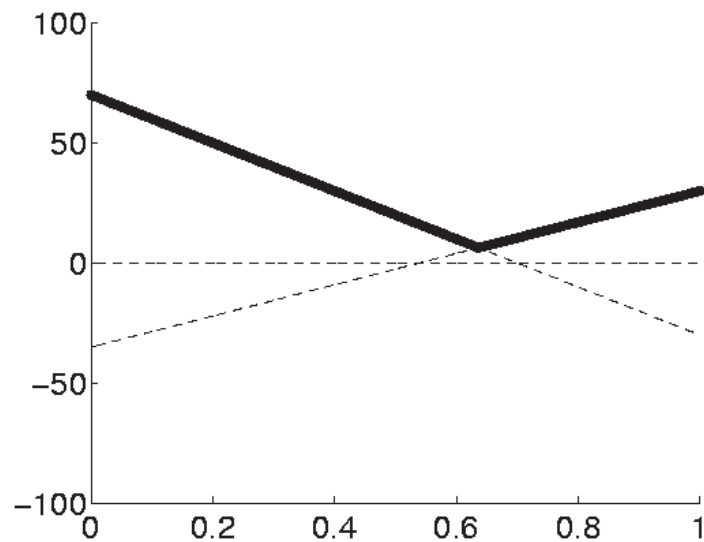
$$\bar{V}_1(b) = \max \left\{ \begin{array}{llll} -70\,p_1 & +30\,(1-p_1) & -30\,p_1 & +70\,(1-p_1) \\ -70\,p_1 & +30\,(1-p_1) & +30\,p_1 & -35\,(1-p_1) \\ +70\,p_1 & -15\,(1-p_1) & -30\,p_1 & +70\,(1-p_1) \\ +70\,p_1 & -15\,(1-p_1) & +30\,p_1 & -35\,(1-p_1) \end{array} \right\}$$

$$= \max \left\{ \begin{array}{ll} -100\,p_1 & +100\,(1-p_1) \\ +40\,p_1 & +55\,(1-p_1) \\ +100\,p_1 & -50\,(1-p_1) \end{array} \right\}$$

# Value Function



$p(z_1) \ V_1(b|z_1)$
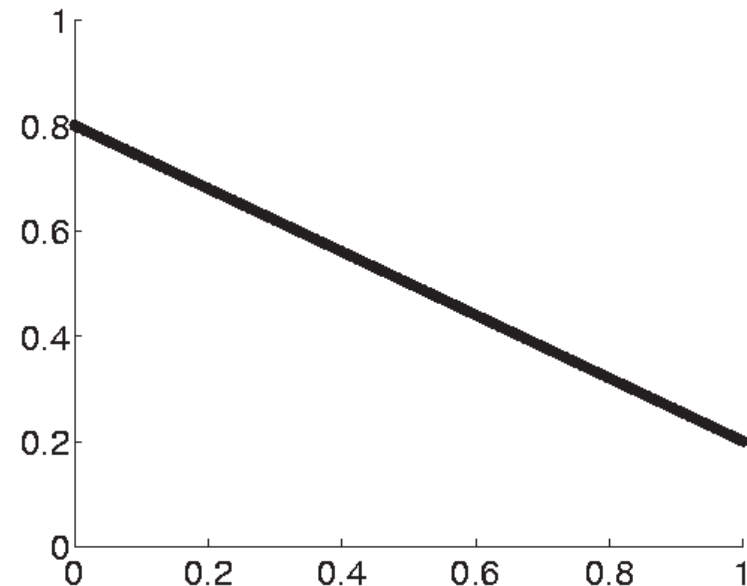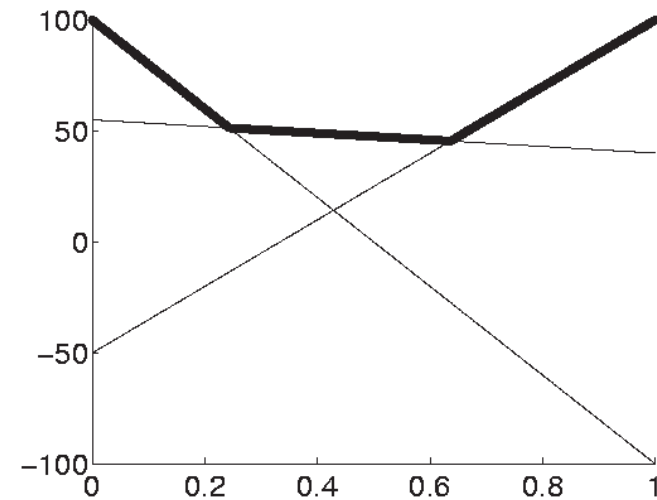
$p(z_2) \ V_2(b|z_2)$

$\bar{V}_1(b)$

74

# 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.

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

# State Transitions (Prediction)

$$
\begin{aligned}
p'_1 &= E_x[p(x_1 \mid x, u_3)] \\
&= \sum_{i=1}^{2} p(x_1 \mid x_i, u_3)p_i \\
&= 0.2p_1 + 0.8(1 - p_1) \\
&= 0.8 - 0.6p_1
\end{aligned}
$$

# Resulting Value Function after executing $u_3$

- Taking the state transitions into account, we finally obtain.

$$\bar{V}_1(b) = \max \begin{cases} -70\,p_1 & +30\,(1-p_1) & -30\,p_1 & +70\,(1-p_1) \\ -70\,p_1 & +30\,(1-p_1) & +30\,p_1 & -35\,(1-p_1) \\ +70\,p_1 & -15\,(1-p_1) & -30\,p_1 & +70\,(1-p_1) \\ +70\,p_1 & -15\,(1-p_1) & +30\,p_1 & -35\,(1-p_1) \end{cases}$$

$$= \max \begin{cases} -100\,p_1 & +100\,(1-p_1) \\ +40\,p_1 & +55\,(1-p_1) \\ +100\,p_1 & -50\,(1-p_1) \end{cases}$$

$$\bar{V}_1(b \mid u_3) = \max \begin{cases} 60\,p_1 & -60\,(1-p_1) \\ 52\,p_1 & +43\,(1-p_1) \\ -20\,p_1 & +70\,(1-p_1) \end{cases}$$
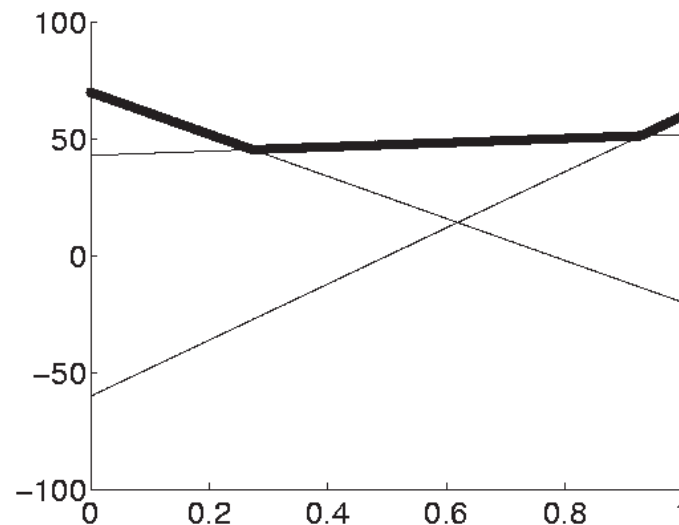
# Value Function after executing $u_3$
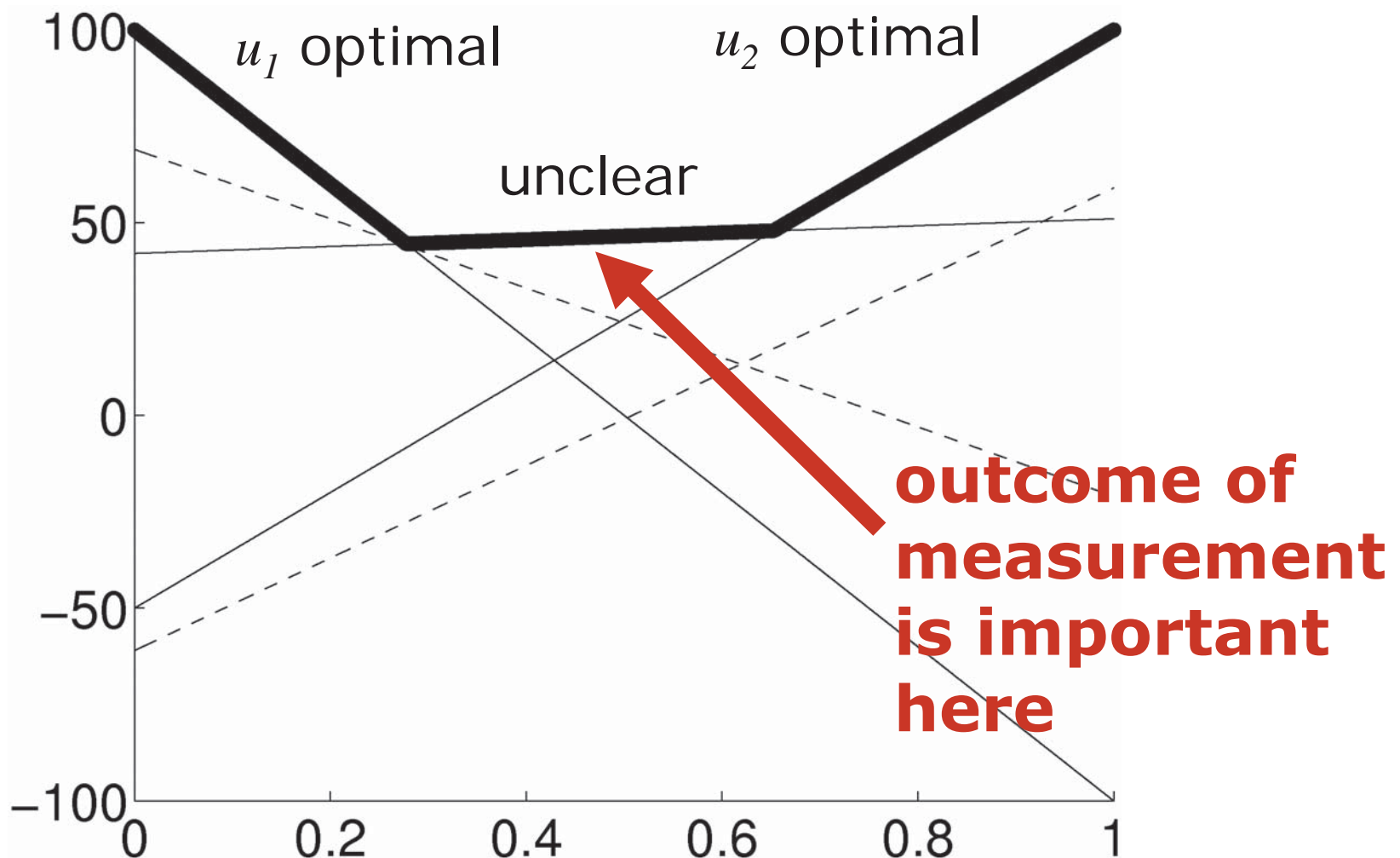
$$\bar{V}_1(b)$$



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

# 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)
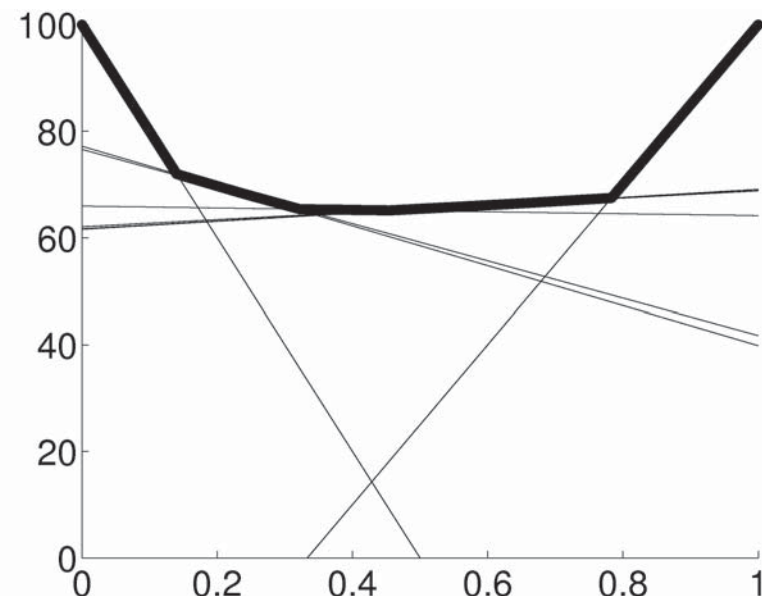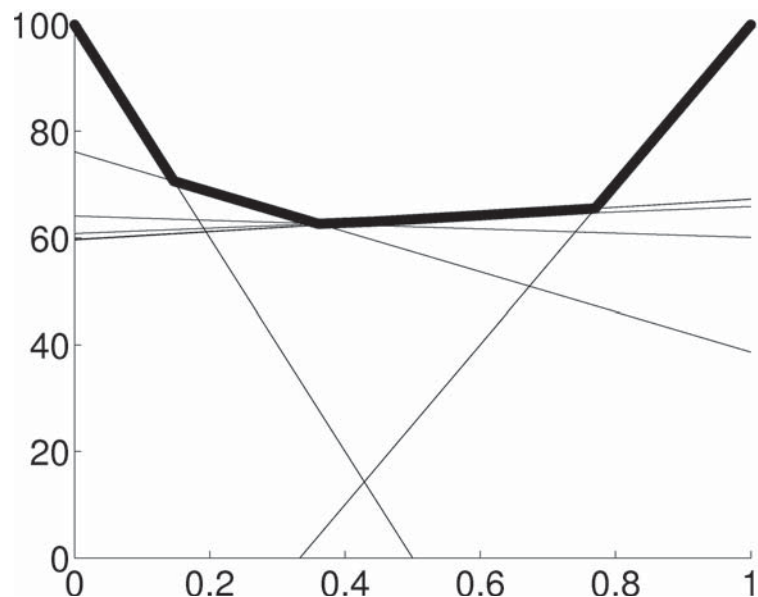
$$\bar{V}_2(b) = \max \left\{ \begin{array}{ll} -100\,p_1 & +100\,(1-p_1) \\ 100\,p_1 & -50\,(1-p_1) \\ 51\,p_1 & +42\,(1-p_1) \end{array} \right\}$$
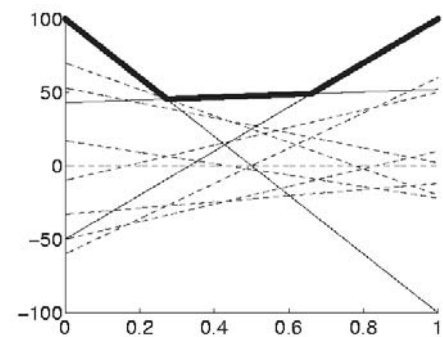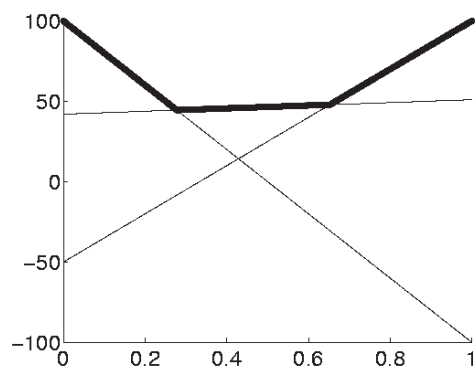
# Graphical Representation of $V_2(b)$



$u_1$ optimal

$u_2$ optimal

unclear

outcome of measurement is important here

# Deep Horizons and Pruning

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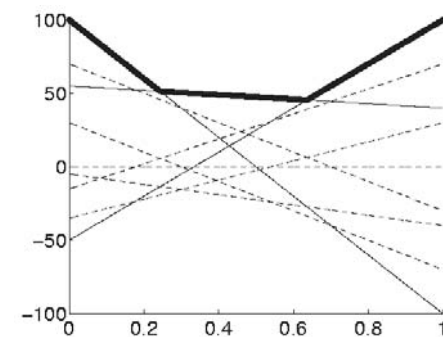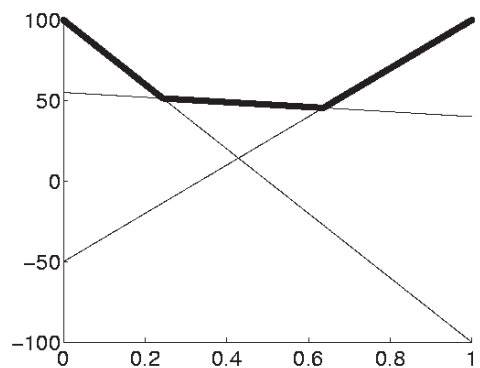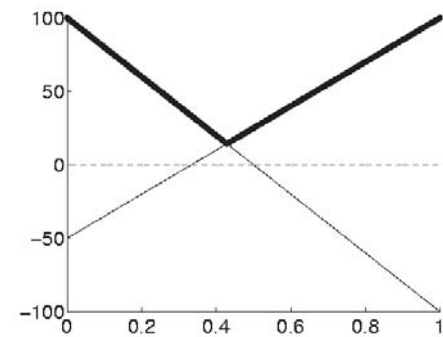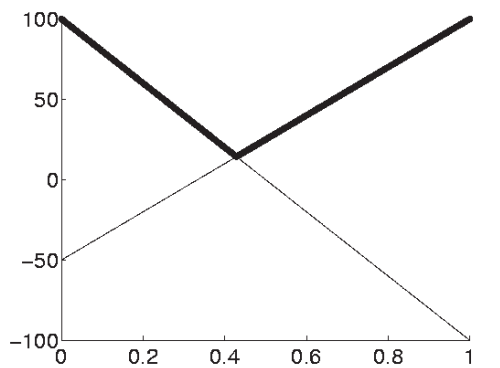- The value functions for T=10 and T=20 are

# Deep Horizons and Pruning



82

- $|S| = 3$
- Hyper-planes
- Finite number of regions over the simplex

- *Sample value function for $|S| = 3$*

- Repeat the process for value functions of 3-horizon,…, and k-horizon POMDP

$$V_t^*(b) = \max_{a \in A} [\sum_i b_i q_i^a + \sum_{i,j,z} b_i p_{ij}^a r_{jz}^a V_{t-1}^*[T(b \mid a, z)]]$$

# Alternate value function interpretation

- **A decision tree**
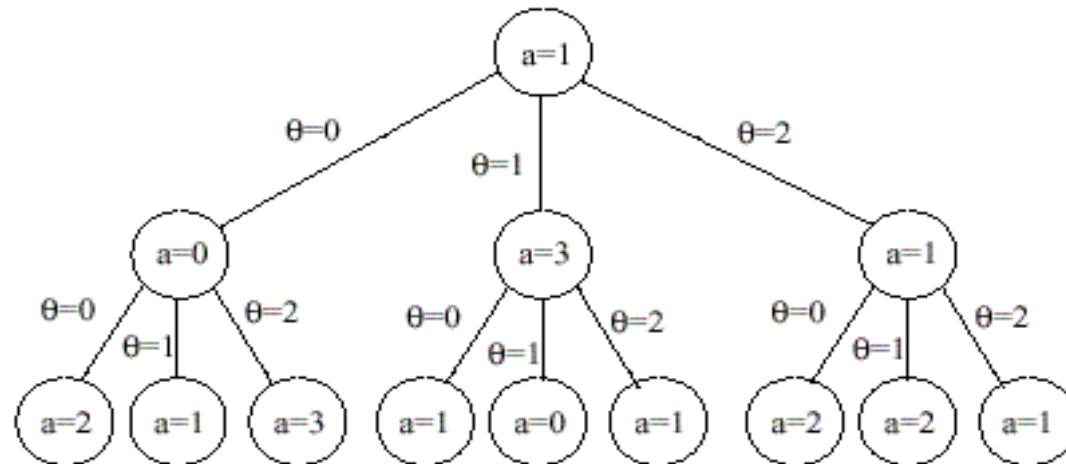  - Nodes represent an action decision
  - Branches represent observation made
- **Too many trees to be generated!**

```
1:          Algorithm POMDP(T):

2:              Υ = (0, . . . , 0)
3:              for τ = 1 to T do
4:                  Υ' = ∅
5:                  for all (u'; v₁ᵏ, . . . , v_Nᵏ) in Υ 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 = Σᵢ₌₁ᴺ vᵢᵏ p(z | xᵢ) p(xᵢ | u, xⱼ)

10:                             endfor
11:                         endfor
12:                     endfor
13:                 endfor
14:                 for all control actions u do
15:                     for all k(1), . . . , k(M) = (1, . . . , 1) to (|Υ|, . . . , |Υ|) do
16:                         for i = 1 to N do

17:                             vᵢ' = γ [ r(xᵢ, u) + Σ_z v_{u,z,i}^{k(z)} ]

18:                         endfor
19:                         add (u; v₁', . . . , v_N') to Υ'
20:                     endfor
21:                 endfor
22:                 optional: prune Υ'
23:                 Υ = Υ'
24:             endfor
25:             return Υ
```

86

# 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.