

# Planning, Learning and Decision Making

Lecture 10. POMDPs (cont.)

# Partially observable MDPs

- **Model** for sequential decision processes
- Described by:
  - State space,  $\mathcal{X}$
  - Action space,  $\mathcal{A}$
  - Observation space,  $\mathcal{Z}$
  - Transition probabilities,  $\{\mathbf{P}_a, a \in \mathcal{A}\}$
  - Observation probabilities,  $\{\mathbf{O}_a, a \in \mathcal{A}\}$
  - Immediate cost function,  $\mathbf{c}$

# Partially observable MDPs

## Key Property: Markov property

The state at instant  $t + 1$  depends only on the state and action at time step  $t$ , i.e.,

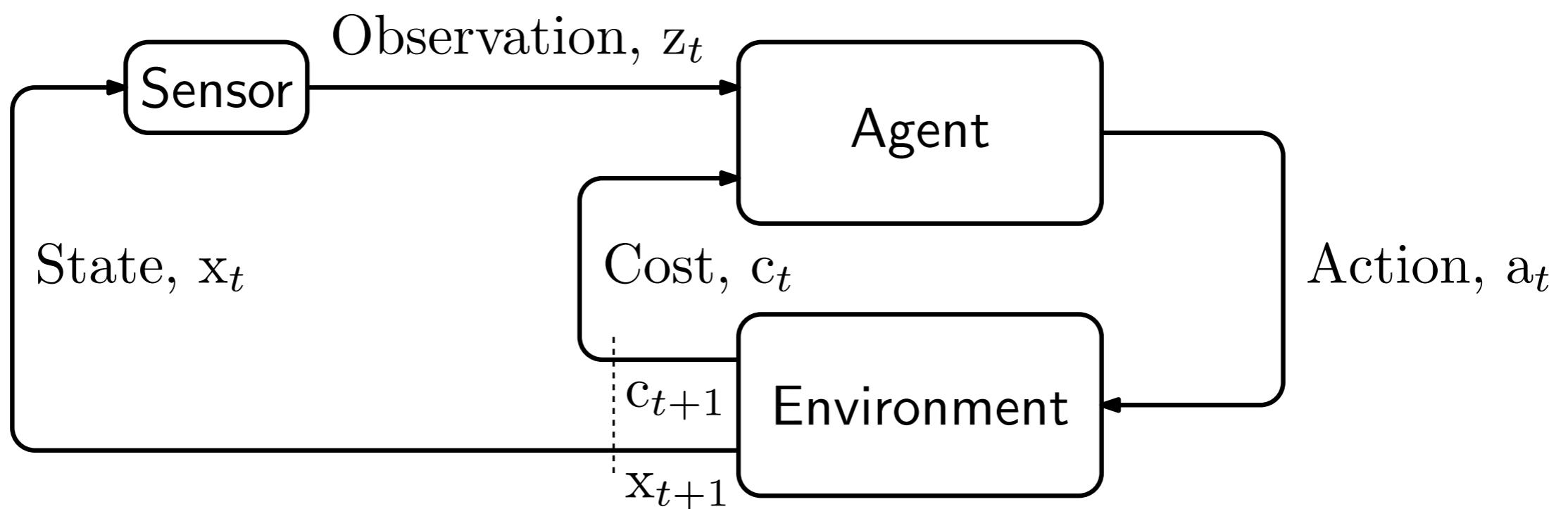
$$\mathbb{P} [x_{t+1} = y \mid x_{0:t} = \mathbf{x}_{0:t}, a_{0:t} = \mathbf{a}_{0:t}] = \mathbb{P} [x_{t+1} = y \mid x_t = x_t, a_t = a_t]$$

## State-dependent observations

The state at instant  $t$  and action at instant  $t - 1$  are enough to predict the observation at instant  $t$ :

$$\begin{aligned} & \mathbb{P} [z_t = z \mid x_{0:t} = \mathbf{x}_{0:t}, a_{0:t-1} = \mathbf{a}_{0:t-1}, z_{0:t-1} = \mathbf{z}_{0:t-1}] = \\ & = \mathbb{P} [z_t = z \mid x_t = x_t, a_{t-1} = a_{t-1}] \end{aligned}$$

# Partially observable MDPs



# Challenges

- Deterministic memoryless policies are not good enough
- Optimal policy may need to keep track of the history...

# A brief recap

- Remember the tiger problem?
  - When escaping a dungeon, you face two doors
  - Behind one of the doors lies your freedom
  - Behind the other door lies a fearsome tiger

# A brief recap

- A POMDP model:

- $\mathcal{X} = \{L, R\}$
- $\mathcal{A} = \{OL, OR, L\}$
- $\mathcal{Z} = \{L, R\}$

- $\dot{\mathbf{P}}_L = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

$$\mathbf{P}_{OL} = \mathbf{P}_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

- $\dot{\mathbf{O}}_L = \begin{bmatrix} 0.85 & 0.15 \\ 0.15 & 0.85 \end{bmatrix}$

$$\mathbf{O}_{OL} = \mathbf{O}_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

- $\dot{\mathbf{C}} = \begin{bmatrix} 1 & 0 & 0.1 \\ 0 & 1 & 0.1 \end{bmatrix}$

# An old trick

- At time  $t = 0$ , you don't know where the tiger is
- You execute "Listen" 2 times
- You observe "Right", "Right"
- How can you use this information?

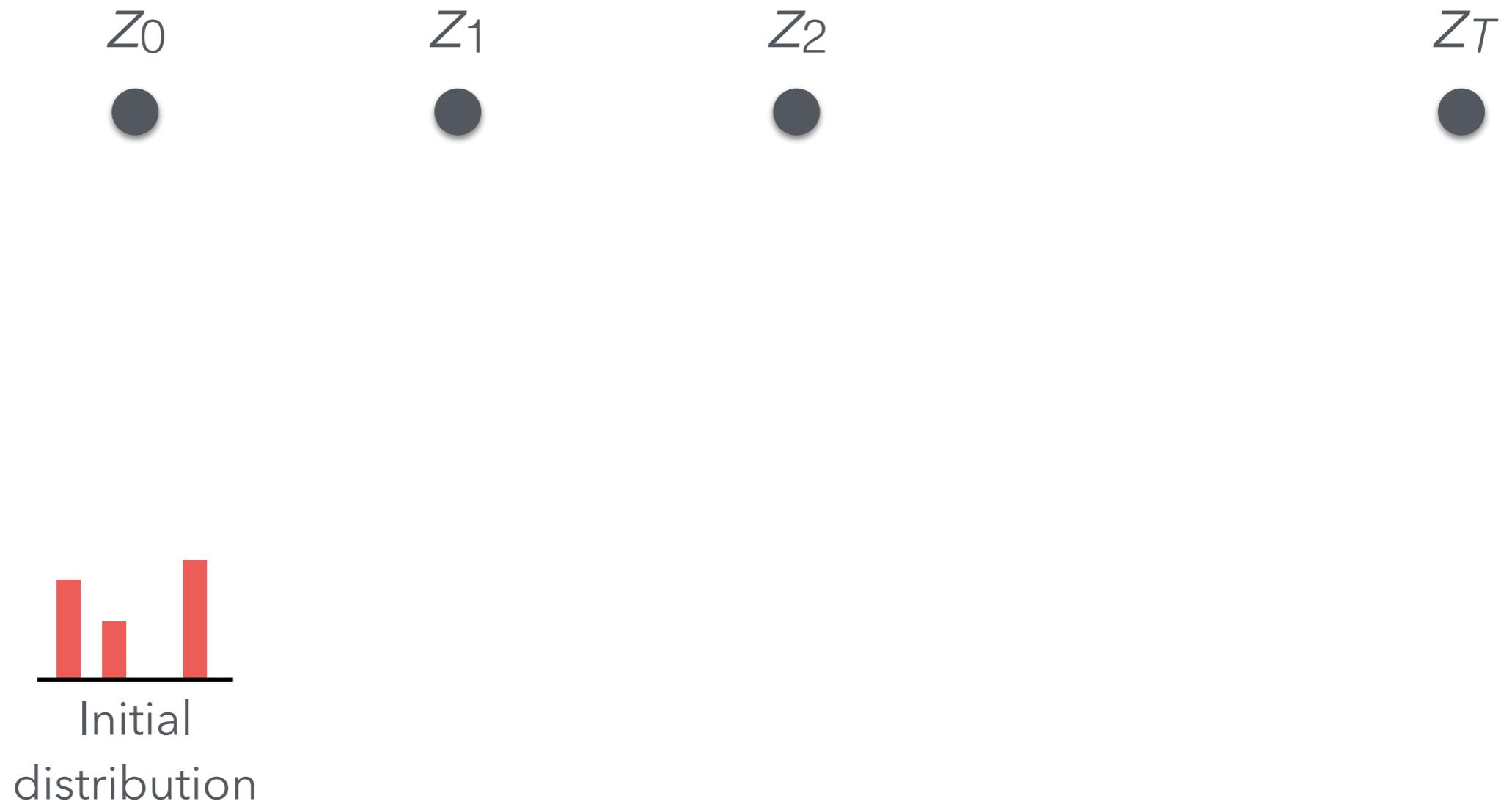
$$\mathbf{P}_{OL} = \mathbf{P}_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$\mathbf{P}_L = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

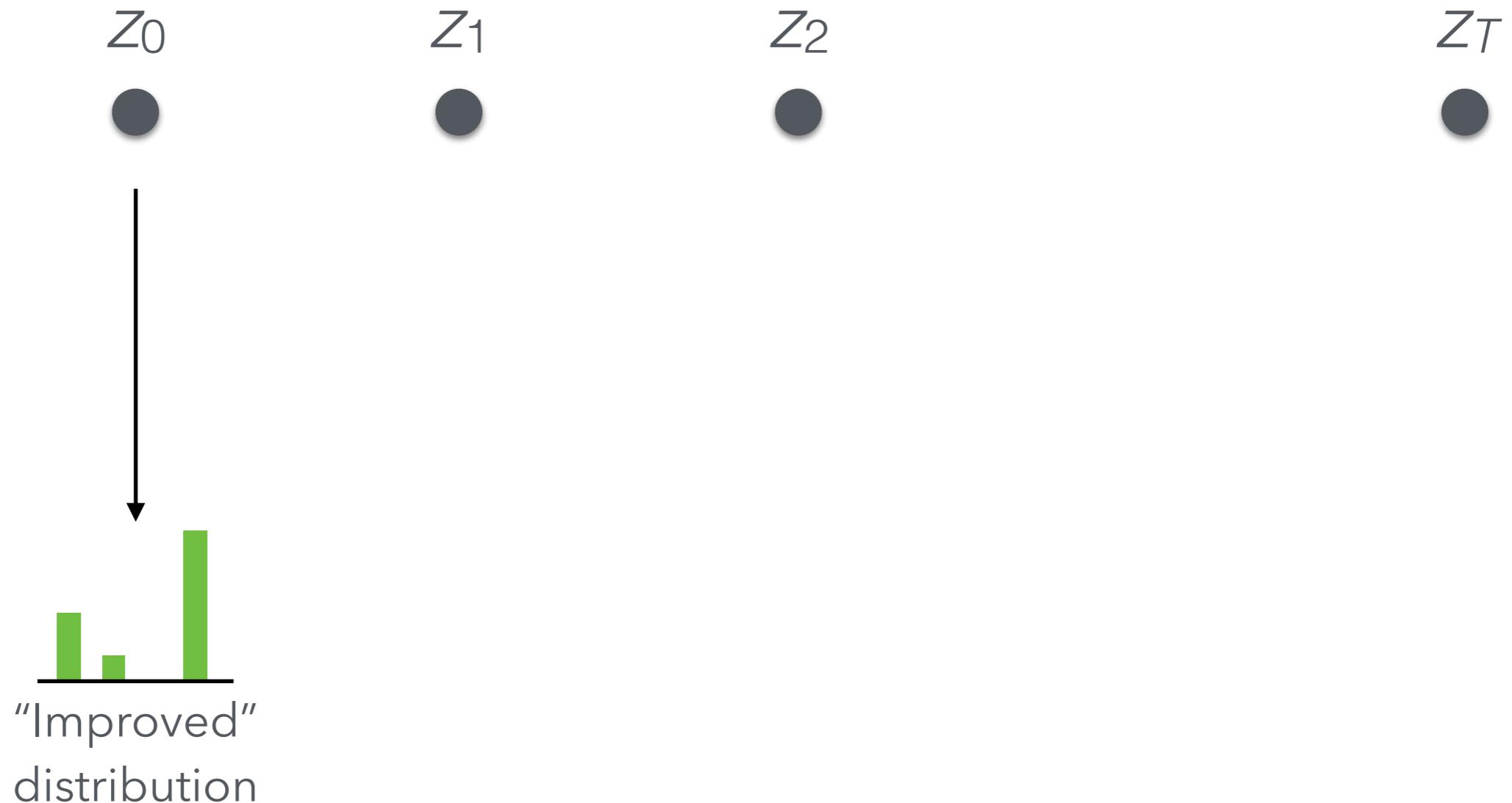
$$\mathbf{O}_{OL} = \mathbf{O}_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$\mathbf{O}_L = \begin{bmatrix} 0.85 & 0.15 \\ 0.15 & 0.85 \end{bmatrix}$$

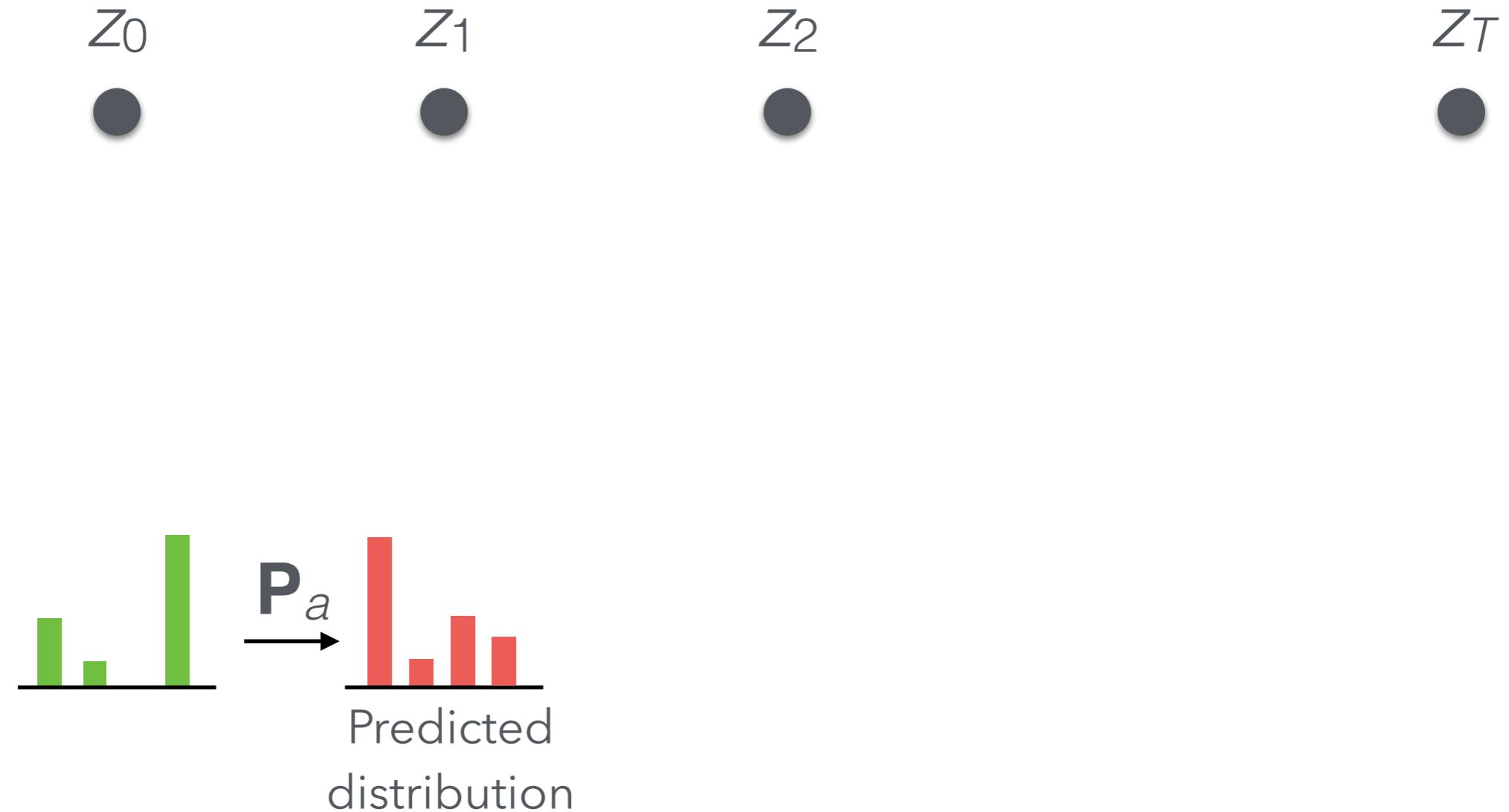
# Remember?



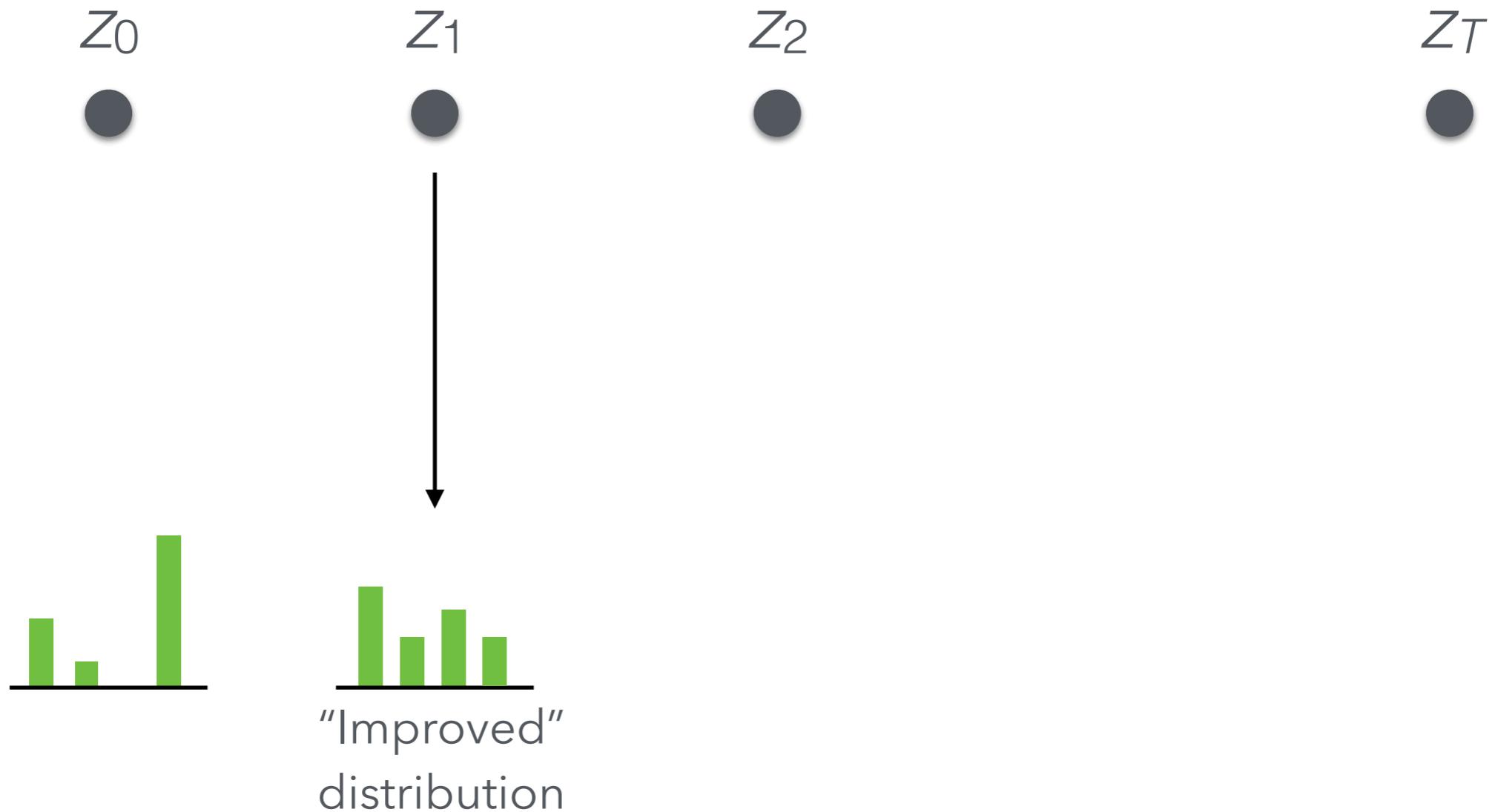
# Remember?



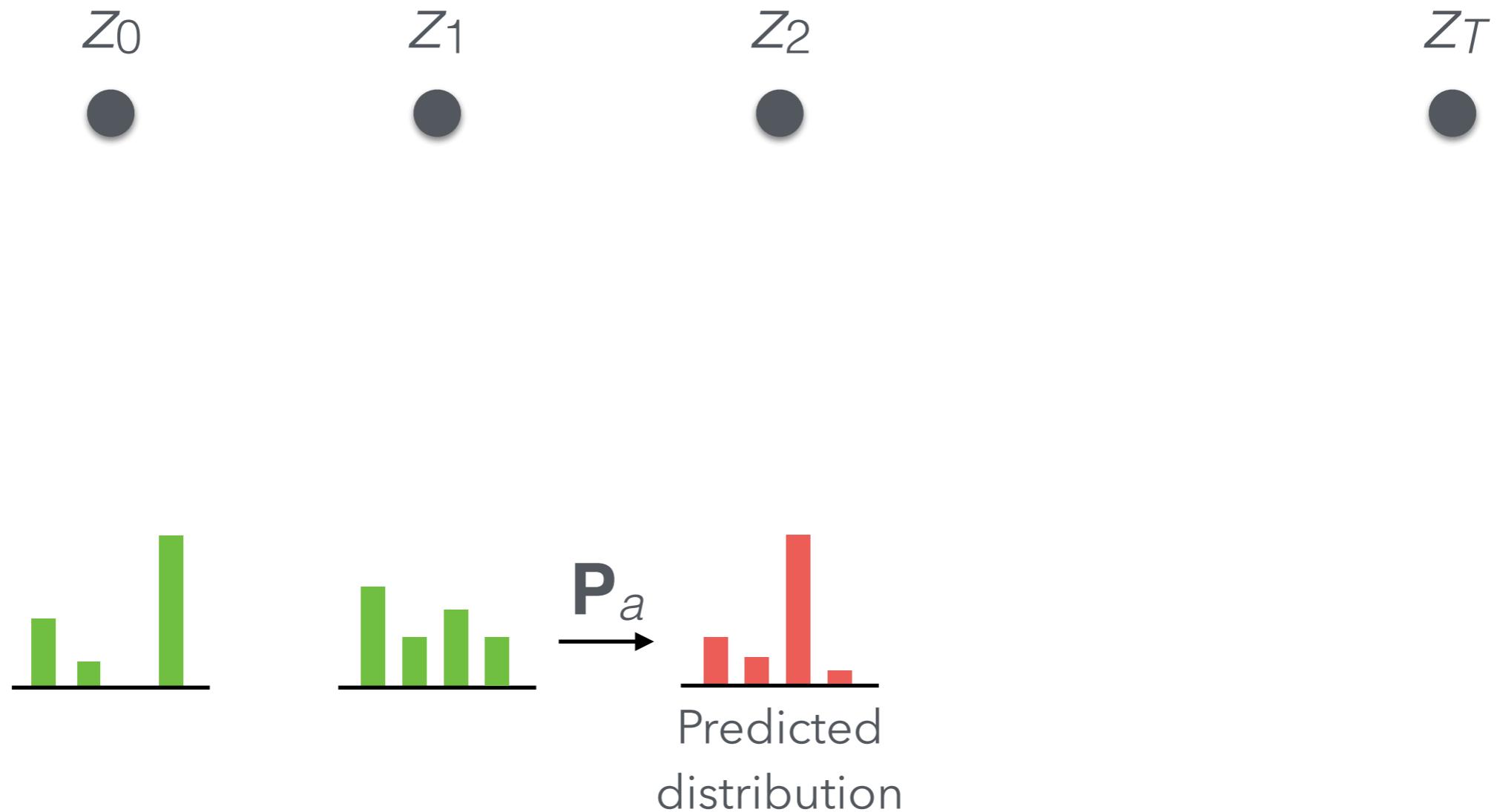
# Remember?



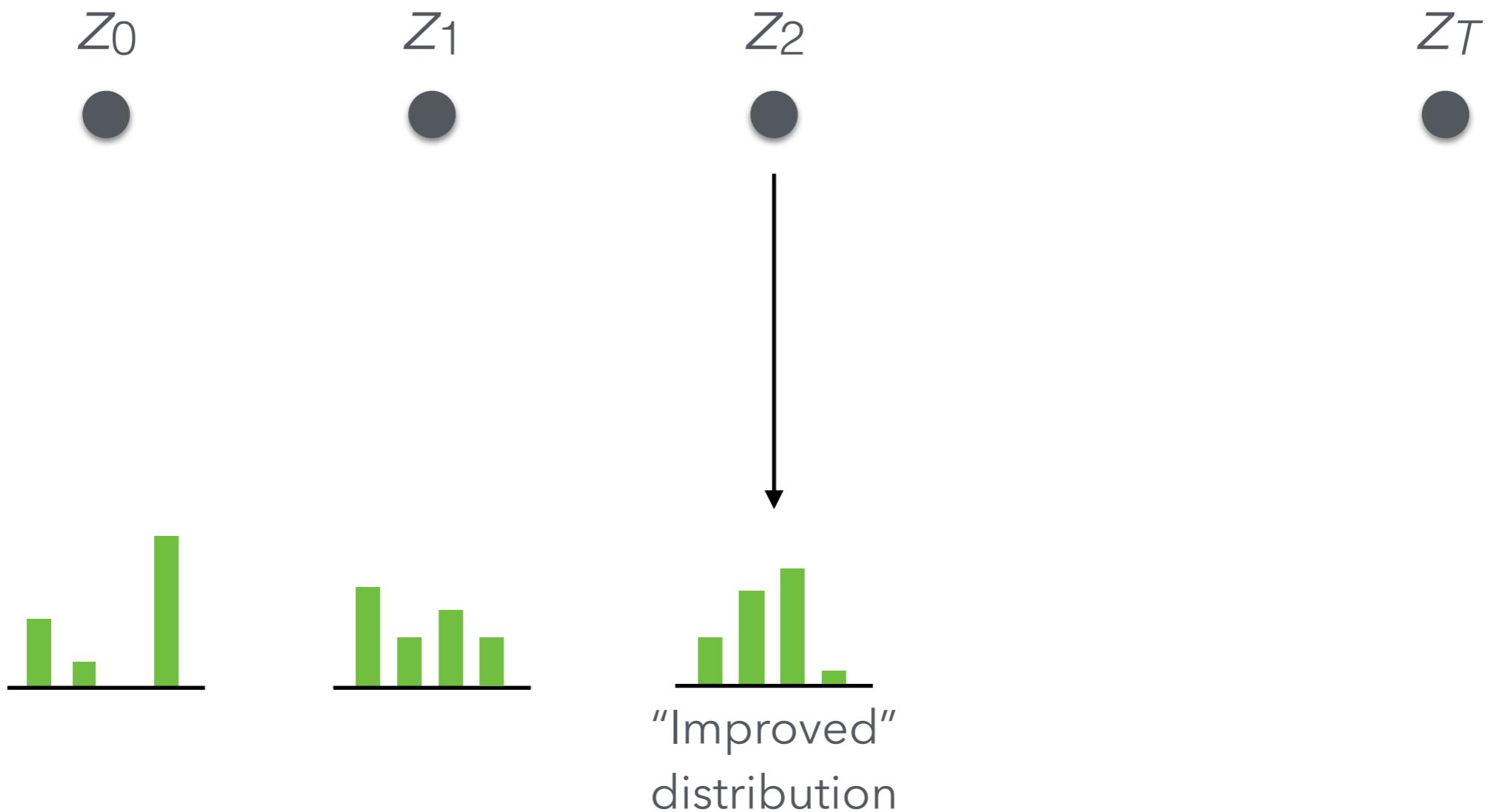
# Remember?



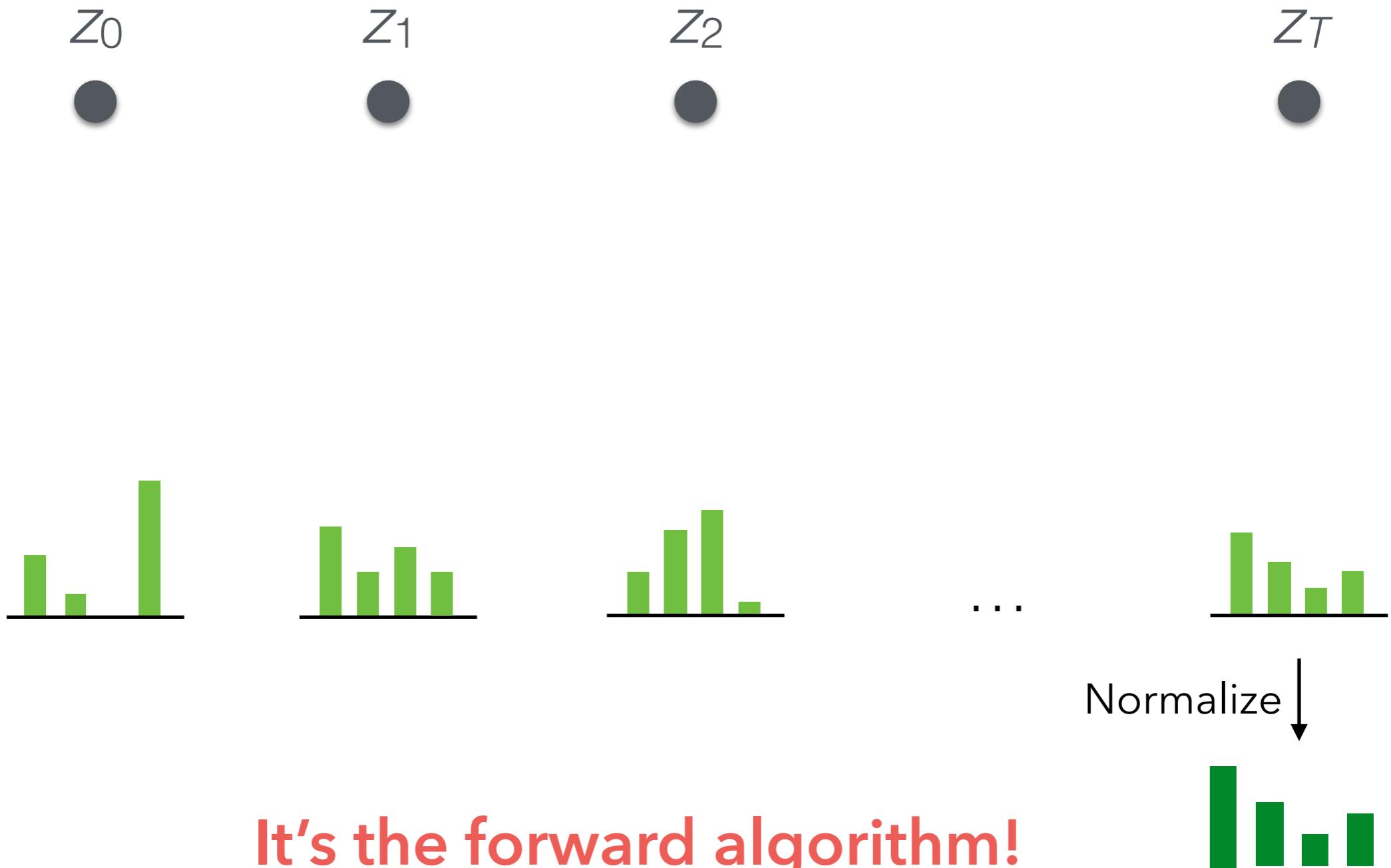
# Remember?



# Remember?



# Remember?



# Let's do this:

- Initial distribution:

$$\alpha_0 = \mu_0 = [0.5 \quad 0.5]$$

$$P_{OL} = P_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$P_L = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$O_{OL} = O_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$O_L = \begin{bmatrix} 0.85 & 0.15 \\ 0.15 & 0.85 \end{bmatrix}$$

# Let's do this:

- Take action “Listen”:

$$\hat{\alpha}_1 = \alpha_0 P_L$$

$$= [0.5 \quad 0.5] \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$= [0.5 \quad 0.5]$$

$$P_{OL} = P_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$P_L = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$O_{OL} = O_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$O_L = \begin{bmatrix} 0.85 & 0.15 \\ 0.15 & 0.85 \end{bmatrix}$$

# Let's do this:

- Consider observation "R":

$$\begin{aligned}\alpha_1 &= \hat{\alpha}_1 \text{diag}(\mathbf{O}_{L,R}) \\ &= [0.5 \quad 0.5] \begin{bmatrix} 0.15 & 0 \\ 0 & 0.85 \end{bmatrix} \\ &= [0.075 \quad 0.425]\end{aligned}$$

$$\mathbf{P}_{OL} = \mathbf{P}_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$\mathbf{P}_L = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mathbf{O}_{OL} = \mathbf{O}_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$\mathbf{O}_L = \begin{bmatrix} 0.85 & 0.15 \\ 0.15 & 0.85 \end{bmatrix}$$

# Let's do this:

- Again, action “Listen”:

$$\hat{\alpha}_2 = \alpha_1 P_L$$

$$= [0.075 \quad 0.425] \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$= [0.075 \quad 0.425]$$

$$P_{OL} = P_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$P_L = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$O_{OL} = O_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$O_L = \begin{bmatrix} 0.85 & 0.15 \\ 0.15 & 0.85 \end{bmatrix}$$

# Let's do this:

- ... and observation "R":

$$\begin{aligned}\alpha_2 &= \hat{\alpha}_2 \text{diag}(\mathbf{O}_{L,R}) \\ &= [0.075 \quad 0.425] \begin{bmatrix} 0.15 & 0 \\ 0 & 0.85 \end{bmatrix} \\ &= [0.011 \quad 0.361]\end{aligned}$$

$$\mathbf{P}_{OL} = \mathbf{P}_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$\mathbf{P}_L = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mathbf{O}_{OL} = \mathbf{O}_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$\mathbf{O}_L = \begin{bmatrix} 0.85 & 0.15 \\ 0.15 & 0.85 \end{bmatrix}$$

# Let's do this:

- Finally, normalize:

$$\begin{aligned}
 \mu_{2|0:2} &= \frac{\alpha_2}{\|\alpha_2\|_1} \\
 &= \frac{1}{0.373} [0.011 \quad 0.361] \\
 &= [0.03 \quad 0.97]
 \end{aligned}$$

$$\mathbf{P}_{OL} = \mathbf{P}_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$\mathbf{P}_L = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mathbf{O}_{OL} = \mathbf{O}_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$\mathbf{O}_L = \begin{bmatrix} 0.85 & 0.15 \\ 0.15 & 0.85 \end{bmatrix}$$

# One extra observation?

- What if we make one extra observation?
  - This time, "L"?

Step     $\mu_{2|0:2} \rightarrow \mu_{2|0:2} P_L$   
 ↓  
 Obs.     $\rightarrow \mu_{2|0:2} P_L \text{diag}(O_{L,L})$   
 ↓  
 Norm.     $\rightarrow \frac{\mu_{2|0:2} P_L \text{diag}(O_{L,L})}{\|\mu_{2|0:2} P_L \text{diag}(O_{L,L})\|_1}$

$$P_{OL} = P_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$P_L = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$O_{OL} = O_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$O_L = \begin{bmatrix} 0.85 & 0.15 \\ 0.15 & 0.85 \end{bmatrix}$$

# One extra observation?

- What if we make one extra observation?
  - This time, "L"?

$$\mu_{3|0:3} = [0.15 \quad 0.85]$$



Probability that tiger is on the right.

$$P_{OL} = P_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$P_L = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$O_{OL} = O_{OR} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$O_L = \begin{bmatrix} 0.85 & 0.15 \\ 0.15 & 0.85 \end{bmatrix}$$

# The belief

- We call the distribution  $\mu_{t|0:t}$  the **belief** at time  $t$ 
  - We will denote it as  $\mathbf{b}_t$
  - $\mathbf{b}_t$  is a probability distribution over  $\mathcal{X}$
  - $\mathbf{b}_t(x)$  is the agent's **belief** that  $x_t = x$ , given all the history

# The belief

- We can update the belief using the previous equation
  - ... after executing action a...
  - ... after making observation z...

$$\mathbf{b}_{t+1} = \frac{\mathbf{b}_t \mathbf{P}_a \text{diag}(\mathbf{O}_{a,z})}{\|\mathbf{b}_t \mathbf{P}_a \text{diag}(\mathbf{O}_{a,z})\|_1}$$

or, component-wise, ...

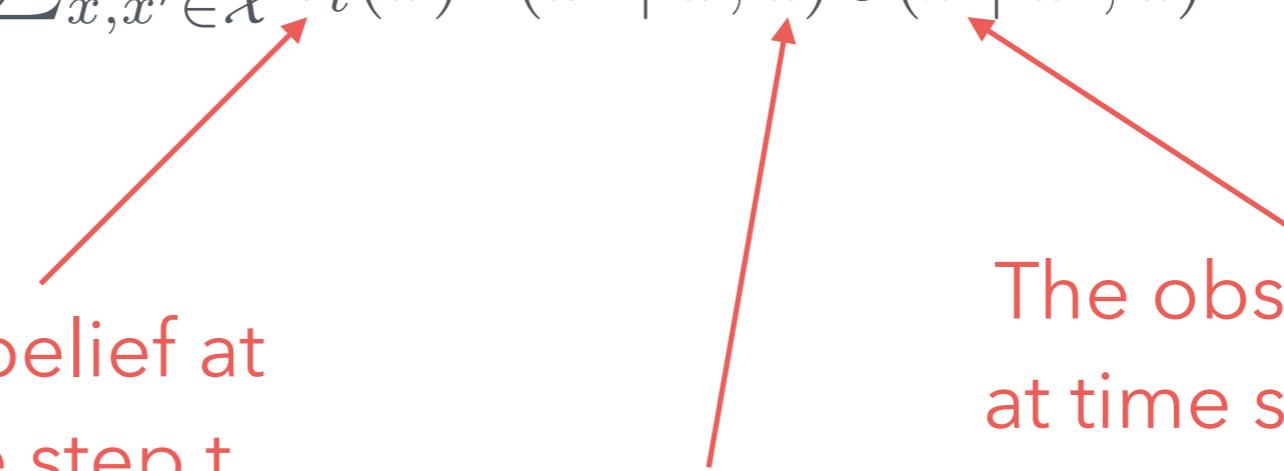
$$\mathbf{b}_{t+1}(y) = \frac{\sum_{x \in \mathcal{X}} b_t(x) \mathbf{P}(y | x, a) \mathbf{O}(z | y, a)}{\sum_{x, x' \in \mathcal{X}} b_t(x) \mathbf{P}(x' | x, a) \mathbf{O}(z | x', a)}$$

Belief  
update

• • •

# The belief

- The belief at time-step  $t + 1$  depends only on:

$$b_{t+1}(y) = \frac{\sum_{x \in \mathcal{X}} b_t(x) P(y | x, a) O(z | y, a)}{\sum_{x, x' \in \mathcal{X}} b_t(x) P(x' | x, a) O(z | x', a)}$$


The belief at time step  $t$

The action at time step  $t$

The observation at time step  $t + 1$

... only quantities from time  $t$

# The belief

- The belief at time-step  $t + 1$  depends only on:

$$b_{t+1}(y) = \frac{\sum_{x \in \mathcal{X}} b_t(x) P(y | x, a) O(z | y, a)}{\sum_{x, x' \in \mathcal{X}} b_t(x) P(x' | x, a) O(z | x', a)}$$

- ... the action at time  $t$
- ... the observation at time  $t$
- ... the belief at time  $t + 1$
- The belief at time  $t$  **summarizes the history** up to time  $t$ !

**The belief is Markov!**

# Partially observable MDPs

- **Model** for sequential decision processes
- Described by:
  - State space,  $\mathcal{X}$
  - Action space,  $\mathcal{A}$
  - Observation space,  $\mathcal{Z}$
  - Transition probabilities,  $\{\mathbf{P}_a, a \in \mathcal{A}\}$
  - Observation probabilities,  $\{\mathbf{O}_a, a \in \mathcal{A}\}$
  - Immediate cost function,  $\mathbf{c}$

# Belief MDPs

- **Model** for sequential decision processes
- Described by:
  - State space, (the space of beliefs)  $\mathcal{B}$
  - Action space,  $\mathcal{A}$
  - Transition probabilities (from the belief update)
  - Immediate cost function,  $c_B$

$$c_B(b, a) = \sum_{x \in \mathcal{X}} b(x)c(x, a)$$

# Optimality

# Cost-to-go function

- We can adapt MDP results to POMDPs through belief-MDPs
- Optimal cost-to-go:

$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \left[ c_B(\mathbf{b}, a) + \gamma \sum_{\mathbf{b}'} \mathbb{P}_B(\mathbf{b}' \mid \mathbf{b}, a) J^*(\mathbf{b}') \right]$$



$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \left[ \sum_{x \in \mathcal{X}} \mathbf{b}(x) \left[ c(x, a) \right] + \gamma \sum_{z \in \mathcal{Z}} \sum_{y \in \mathcal{X}} \mathbb{P}(y \mid x, a) \mathbb{O}(z \mid y, a) J^*(\mathbf{b}'_{z,a}) \right]$$

$c_B(\mathbf{b}, a)$

# Cost-to-go function

- We can adapt MDP results to POMDPs through belief-MDPs
- Optimal cost-to-go:

$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \left[ c_B(\mathbf{b}, a) + \gamma \sum_{\mathbf{b}'} \mathbb{P}_B(\mathbf{b}' \mid \mathbf{b}, a) J^*(\mathbf{b}') \right]$$

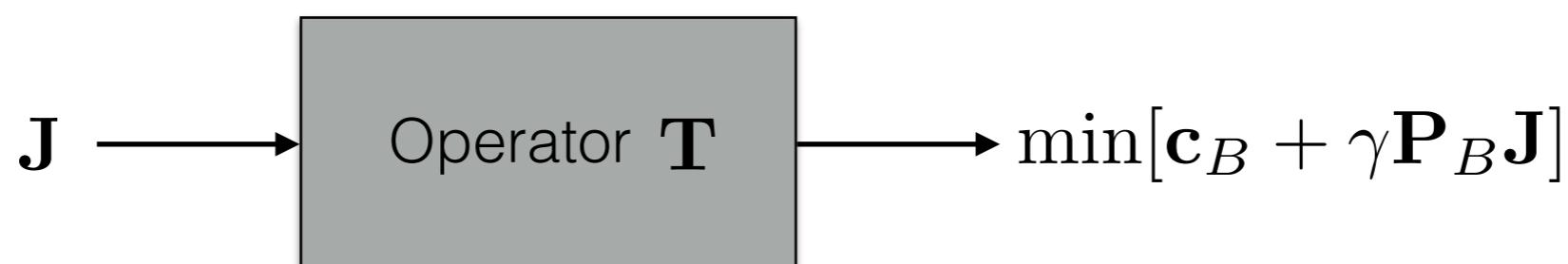
$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \left[ \sum_{x \in \mathcal{X}} \mathbf{b}(x) \left[ c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \sum_{y \in \mathcal{X}} \mathbb{P}(y \mid x, a) \mathbb{O}(z \mid y, a) J^*(\mathbf{b}'_{z,a}) \right] \right]$$

$\mathbb{P}_B(\mathbf{b}' \mid \mathbf{b}, a)$

# Value iteration for POMDPs

- We can use the recursive expression:

$$J^*(\mathbf{b}) = \boxed{\min_{a \in \mathcal{A}} \left[ c_B(\mathbf{b}, a) + \gamma \sum_{\mathbf{b}'} \mathsf{P}_B(\mathbf{b}' \mid \mathbf{b}, a) J^*(\mathbf{b}') \right]}$$



# ... however...

$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \left[ c_B(\mathbf{b}, a) + \gamma \sum_{\mathbf{b}'} P_B(\mathbf{b}' \mid \mathbf{b}, a) J^*(\mathbf{b}') \right]$$

How do we  
represent this?

# Representing $J^*$

- $J^*$  is a function in  $\text{IR}^n$  (belief)
- We cannot represent it explicitly
- Therefore,

---

**Require:** MDP  $\mathcal{M} = (\mathcal{X}, \mathcal{A}, \{\mathbf{P}_x\}_{x \in \mathcal{X}}, \mathbf{R}, \gamma)$ ,  $\gamma \in (0, 1)$ ,  $\varepsilon > 0$ ;

- 1: Initialize  $k = 0$ ,  $J^{(0)} \equiv 0$
- 2: **repeat**
- 3:      $J^{(k+1)} \leftarrow \mathbf{T}J^{(k)}$
- 4:      $k \leftarrow k + 1$
- 5: **until**  $\|J^{(k+1)} - J^{(k)}\|_{\infty} < \varepsilon$ .
- 6: Compute  $\pi^*$  using (2.21)
- 7: **return**  $\pi^*$

---





# The tiger problem

# What is the belief space?

- The POMDP can be in two possible states:
  - Tiger left
  - Tiger right
- Belief  **$b$**  is a vector

$$\mathbf{b} = [\mathbb{P} [\text{Tiger left}] \quad \mathbb{P} [\text{Tiger right}]]$$

# What is the belief space?

- The POMDP can be in two possible states:
  - Tiger left
  - Tiger right
- Belief **b** is a vector

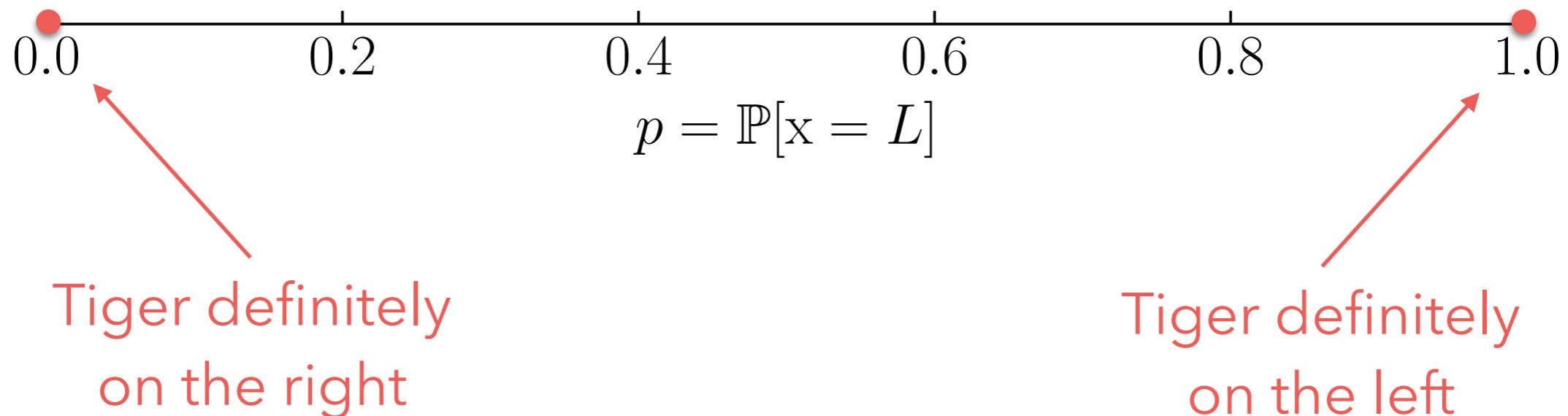
$$\mathbf{b} = [\mathbb{P} [\text{Tiger left}] \quad 1 - \mathbb{P} [\text{Tiger left}]]$$



Just a single  
number

# What is the belief space?

- We can represent it as a number in  $[0, 1]$ :

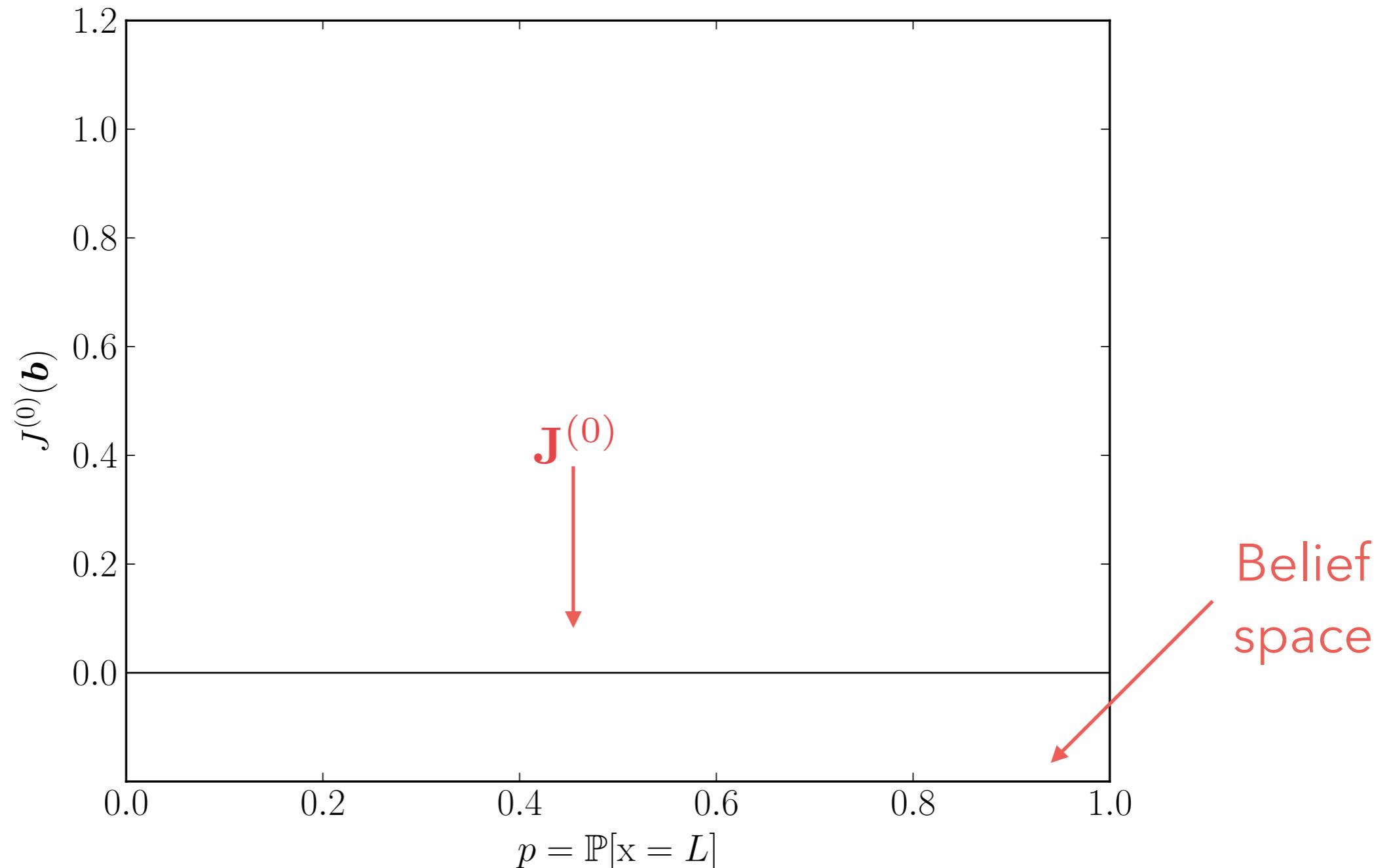


# Value iteration

- Let's try running VI on the tiger problem
  - We start with

$$J^{(0)} \equiv 0$$

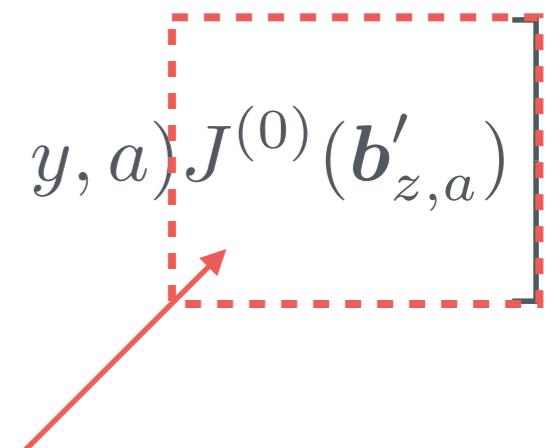
# Value iteration



# Value iteration

- Let's try running VI on the tiger problem
  - Iteration 1:

$$J^{(1)}(\mathbf{b}) = \min_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) \left[ c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \sum_{y \in \mathcal{X}} P(y | x, a) O(z | y, a) J^{(0)}(\mathbf{b}'_{z,a}) \right]$$



This is zero!

# Value iteration

- Let's try running VI on the tiger problem
  - Iteration 1:

$$J^{(1)}(\mathbf{b}) = \min_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) c(x, a)$$

# Value iteration

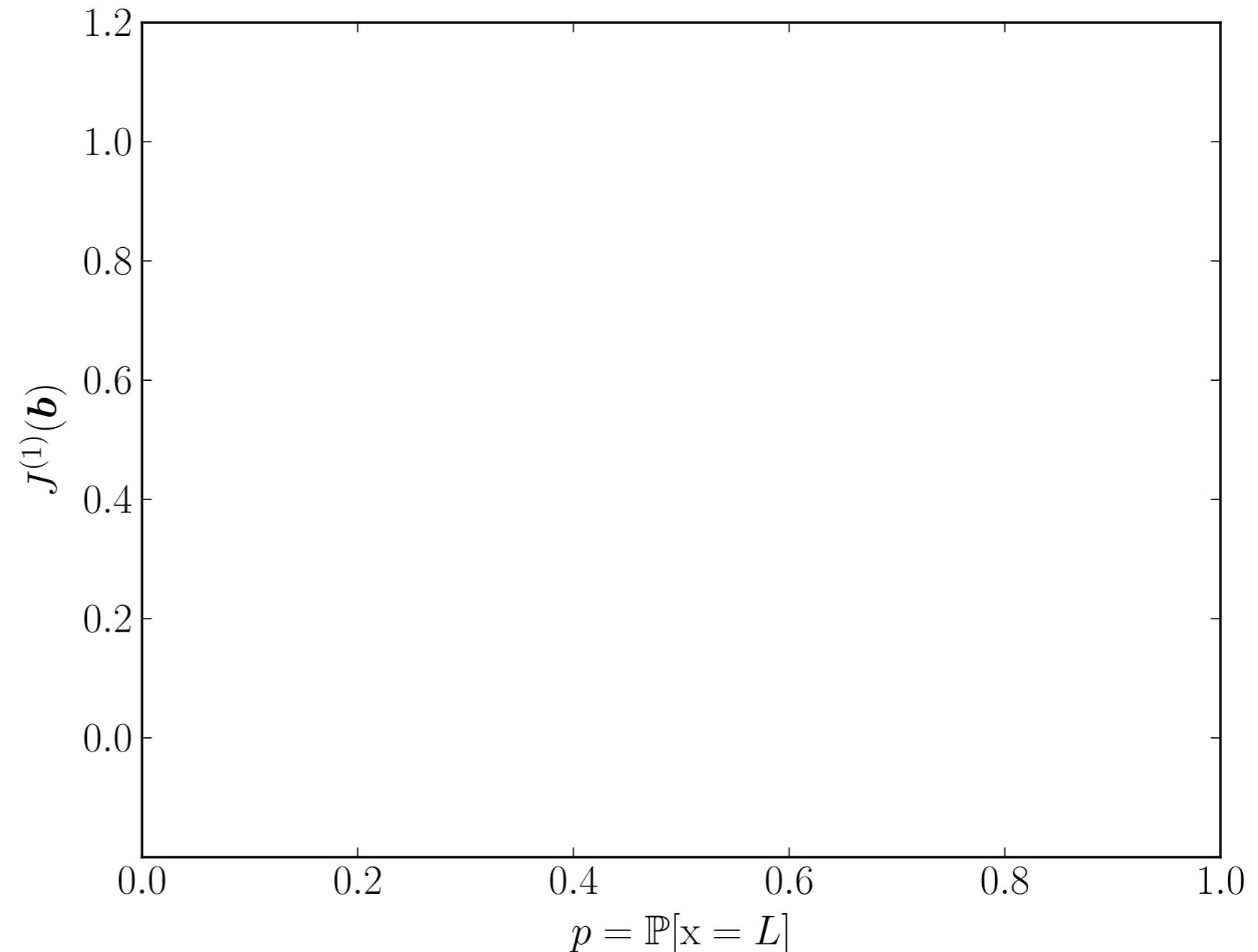
- Let's try running VI on the tiger problem
  - Iteration 1:

$$J^{(1)}(\mathbf{b}) = \min_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) c(x, a)$$

# Value iteration

$$C = \begin{bmatrix} 1 & \boxed{0} & 0.1 \\ 0 & \boxed{1} & 0.1 \end{bmatrix}$$

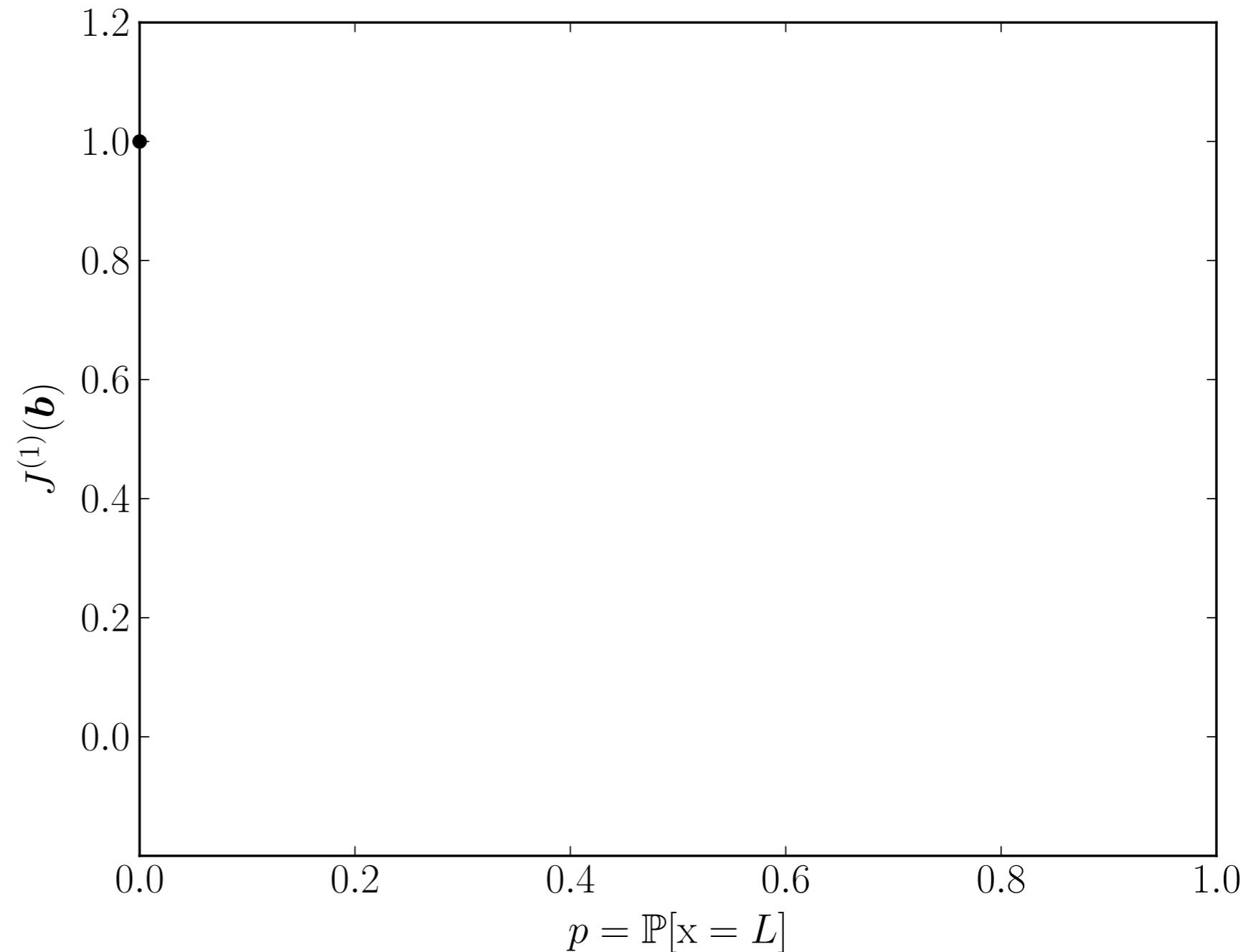
Action OR



# Value iteration

$$C = \begin{bmatrix} 1 & \boxed{0} & 0.1 \\ 0 & \boxed{1} & 0.1 \end{bmatrix}$$

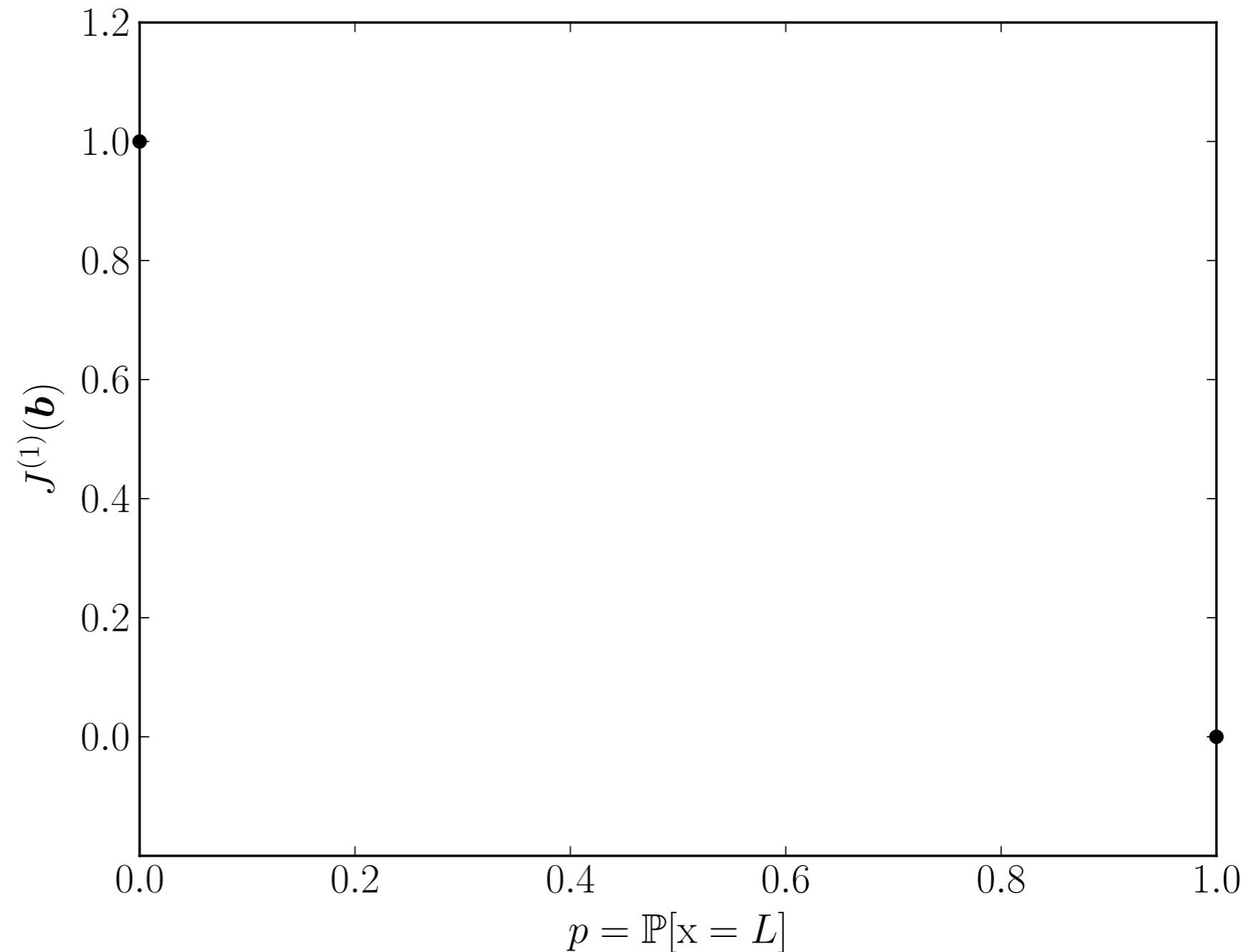
Action OR



# Value iteration

$$C = \begin{bmatrix} 1 & \boxed{0} & 0.1 \\ 0 & \boxed{1} & 0.1 \end{bmatrix}$$

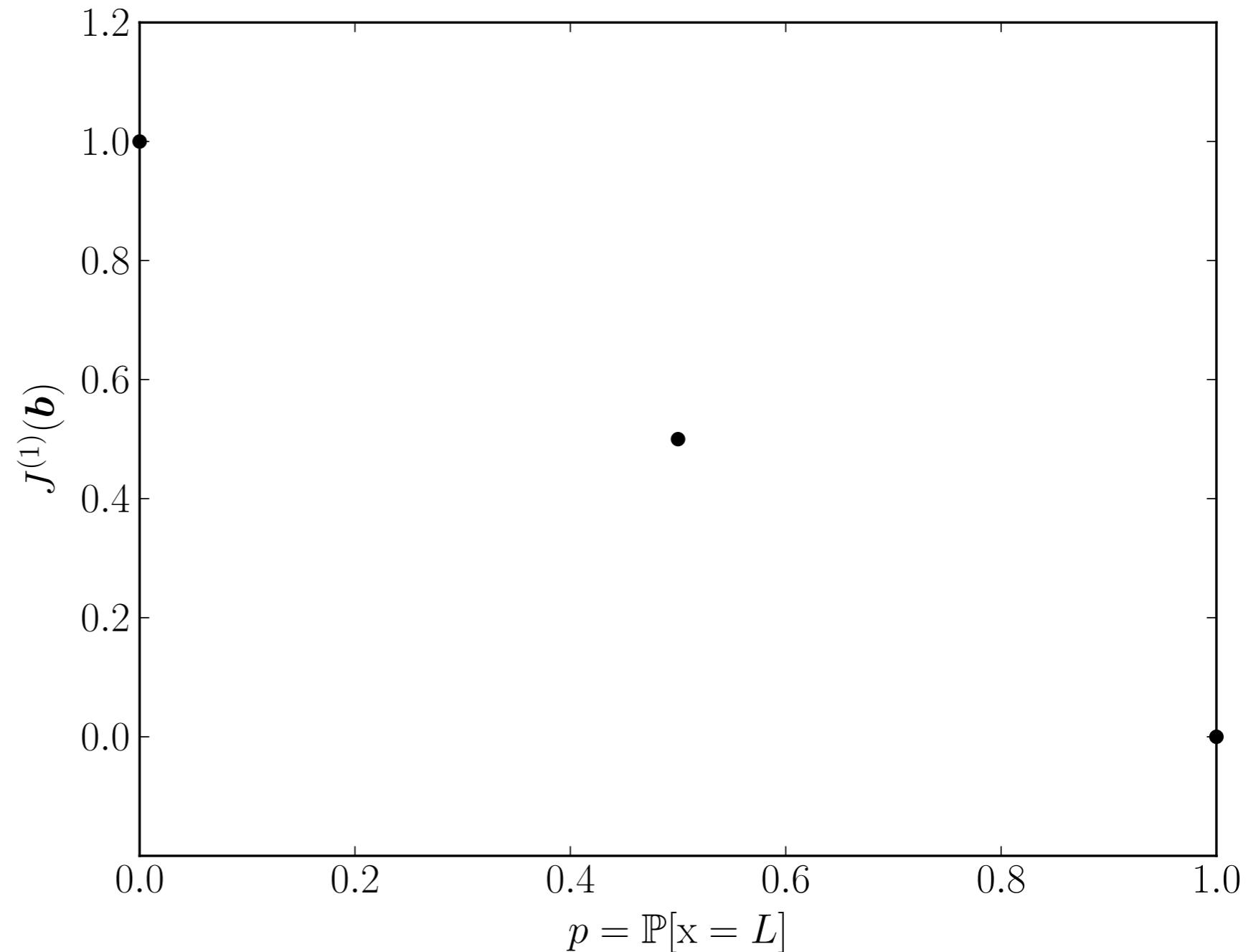
Action OR



# Value iteration

$$C = \begin{bmatrix} 1 & \boxed{0} & 0.1 \\ 0 & \boxed{1} & 0.1 \end{bmatrix}$$

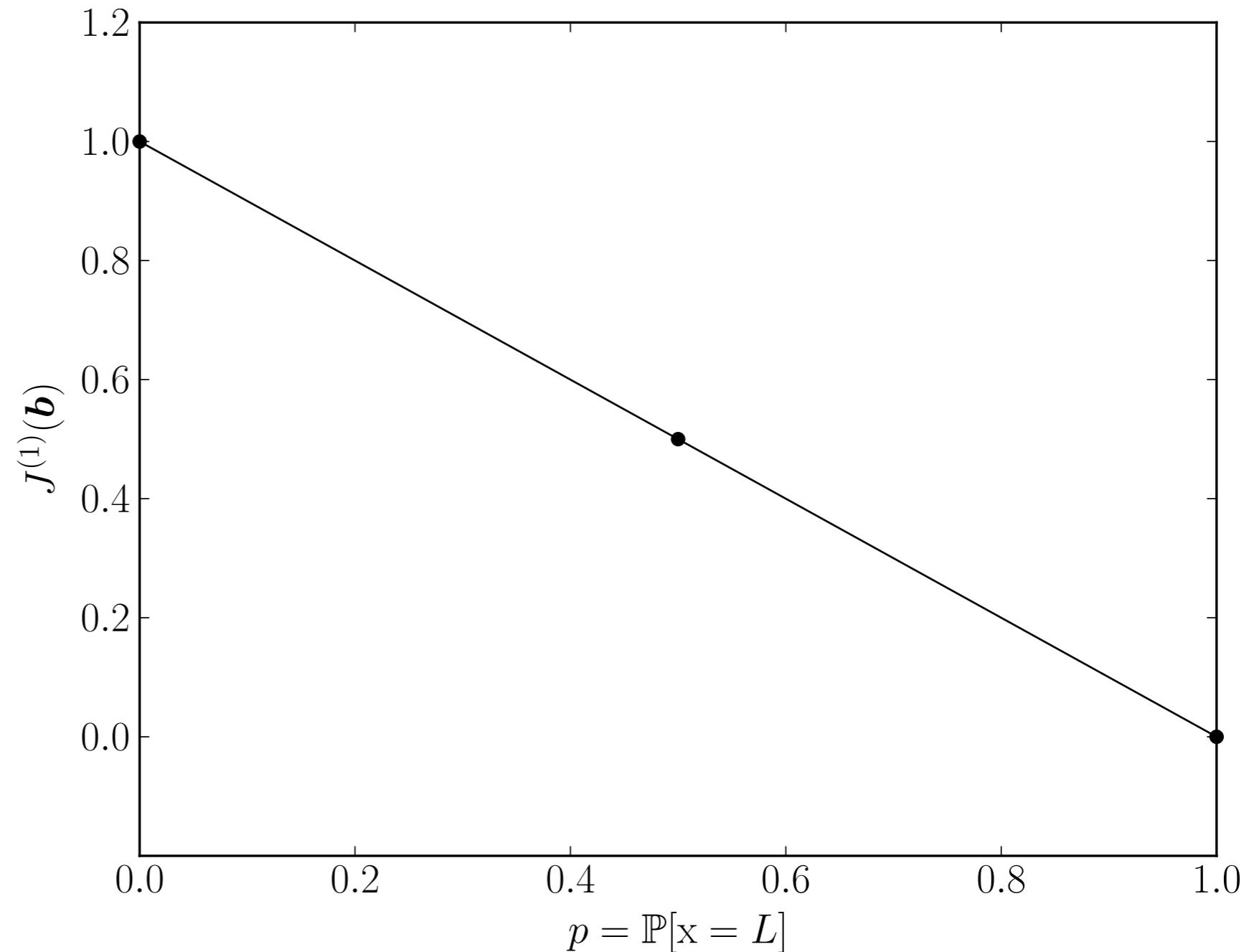
Action OR



# Value iteration

$$C = \begin{bmatrix} 1 & \boxed{0} & 0.1 \\ 0 & \boxed{1} & 0.1 \end{bmatrix}$$

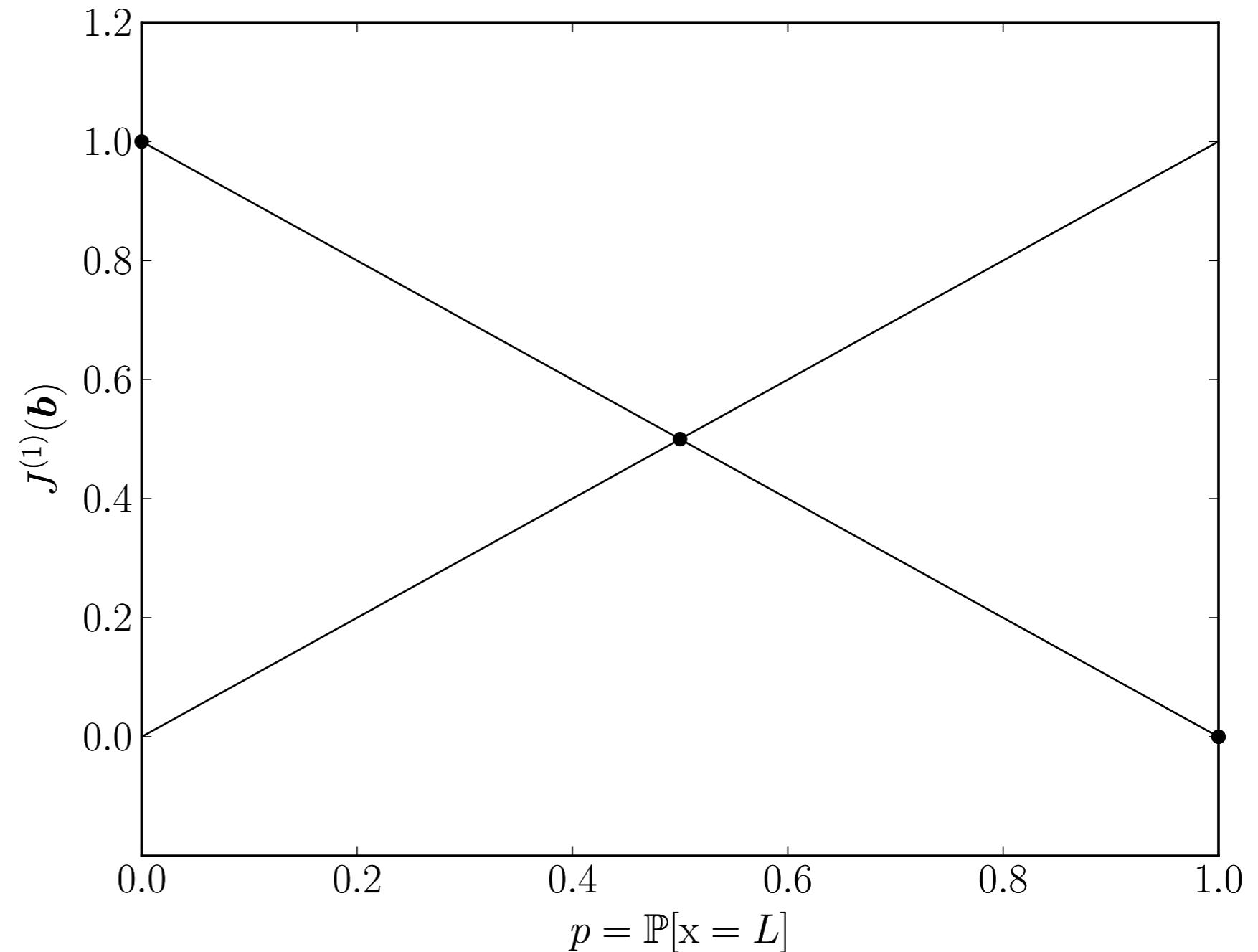
Action OR



# Value iteration

$$C = \begin{bmatrix} 1 & 0 & 0.1 \\ 0 & 1 & 0.1 \end{bmatrix}$$

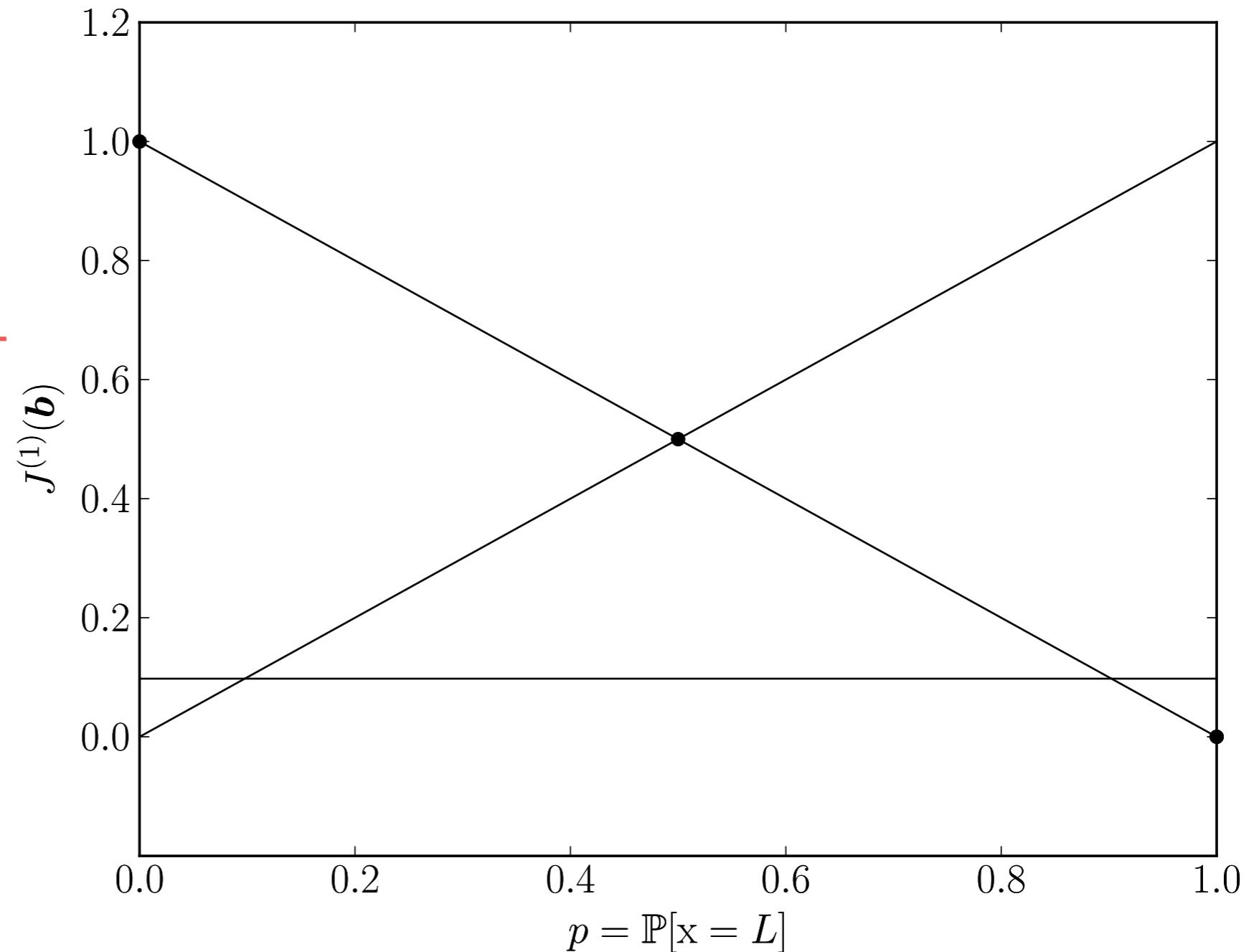
Action OL



# Value iteration

$$C = \begin{bmatrix} 1 & 0 & 0.1 \\ 0 & 1 & 0.1 \end{bmatrix}$$

Action L



# Value iteration

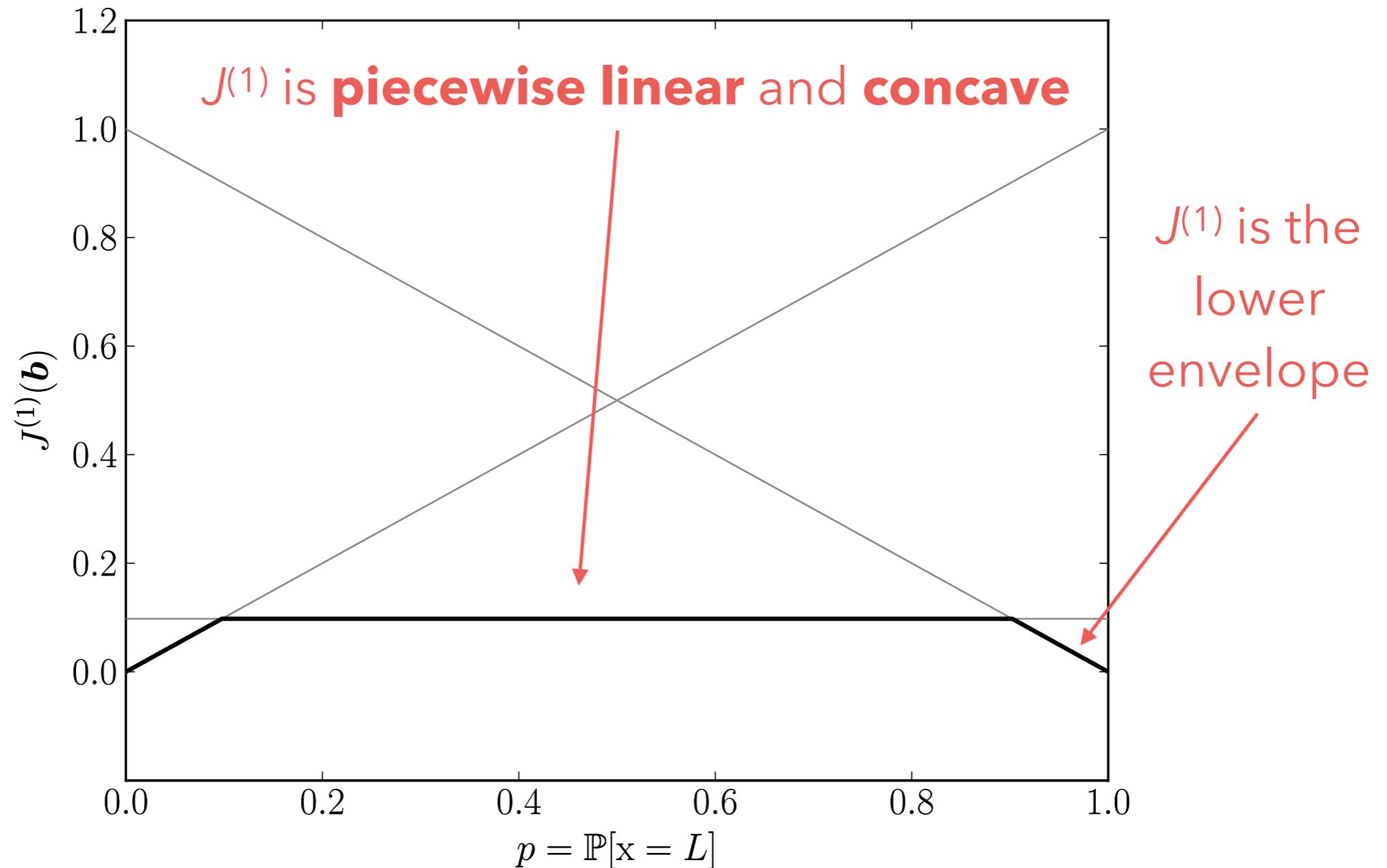
- Iteration 1:

$$J^{(1)}(\mathbf{b}) = \min_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) c(x, a)$$

- We just computed, for every belief  $\mathbf{b}$  and action  $a$ ,

$$\sum_{x \in \mathcal{X}} \mathbf{b}(x) c(x, a)$$

# Value iteration

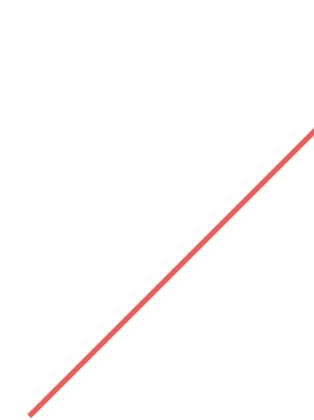


# Representing $J(k)$

- The cost-to-go at each iteration of VI is always PWLC
  - Can always be written in the form

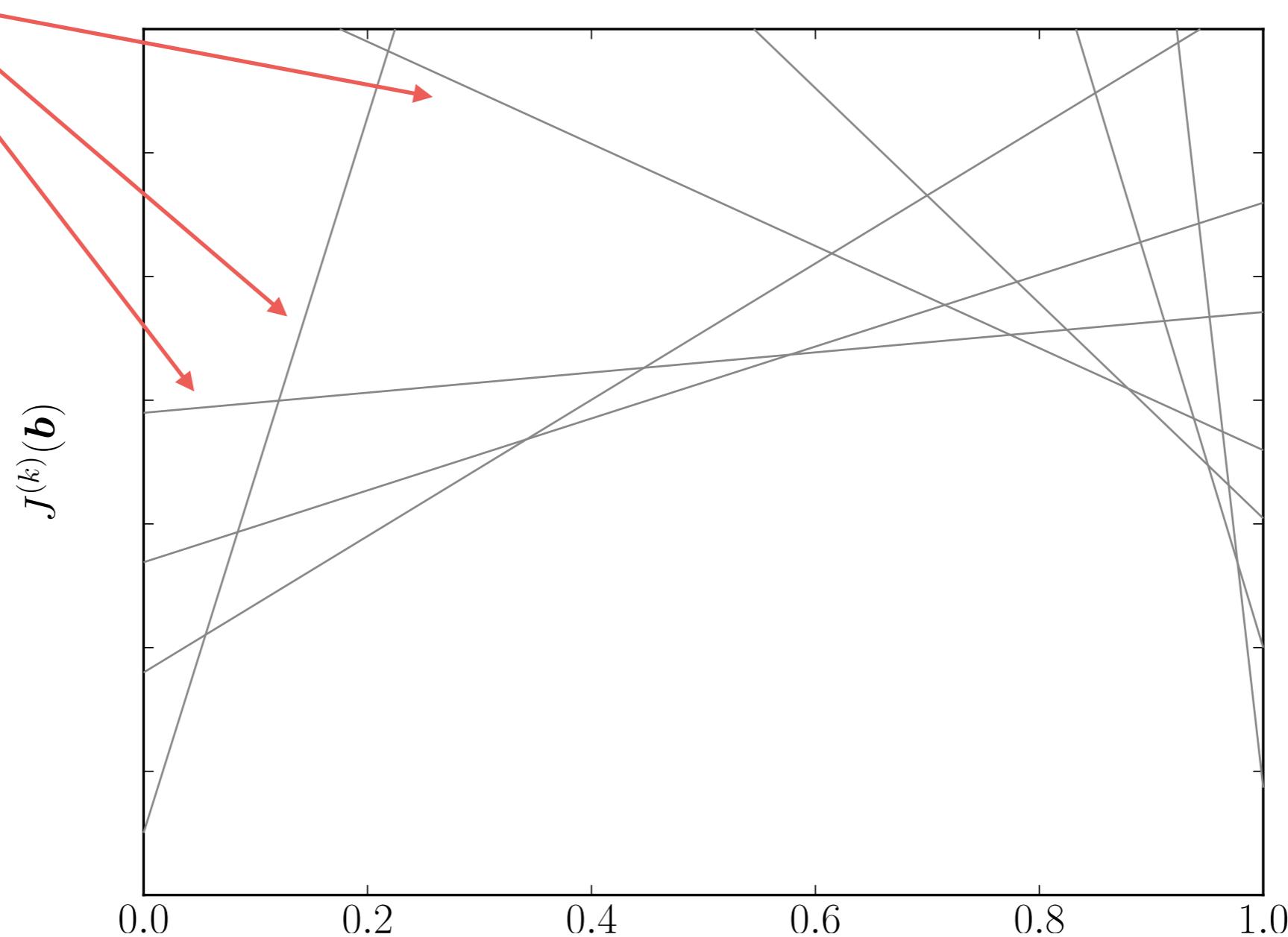
$$\begin{aligned} J^{(k)}(\mathbf{b}) &= \min_{\alpha \in \Gamma} \mathbf{b} \cdot \alpha \\ &= \min_{\alpha \in \Gamma} \sum_{x \in \mathcal{X}} \mathbf{b}(x) \alpha(x) \end{aligned}$$

Set of vectors  
used in the  
representation



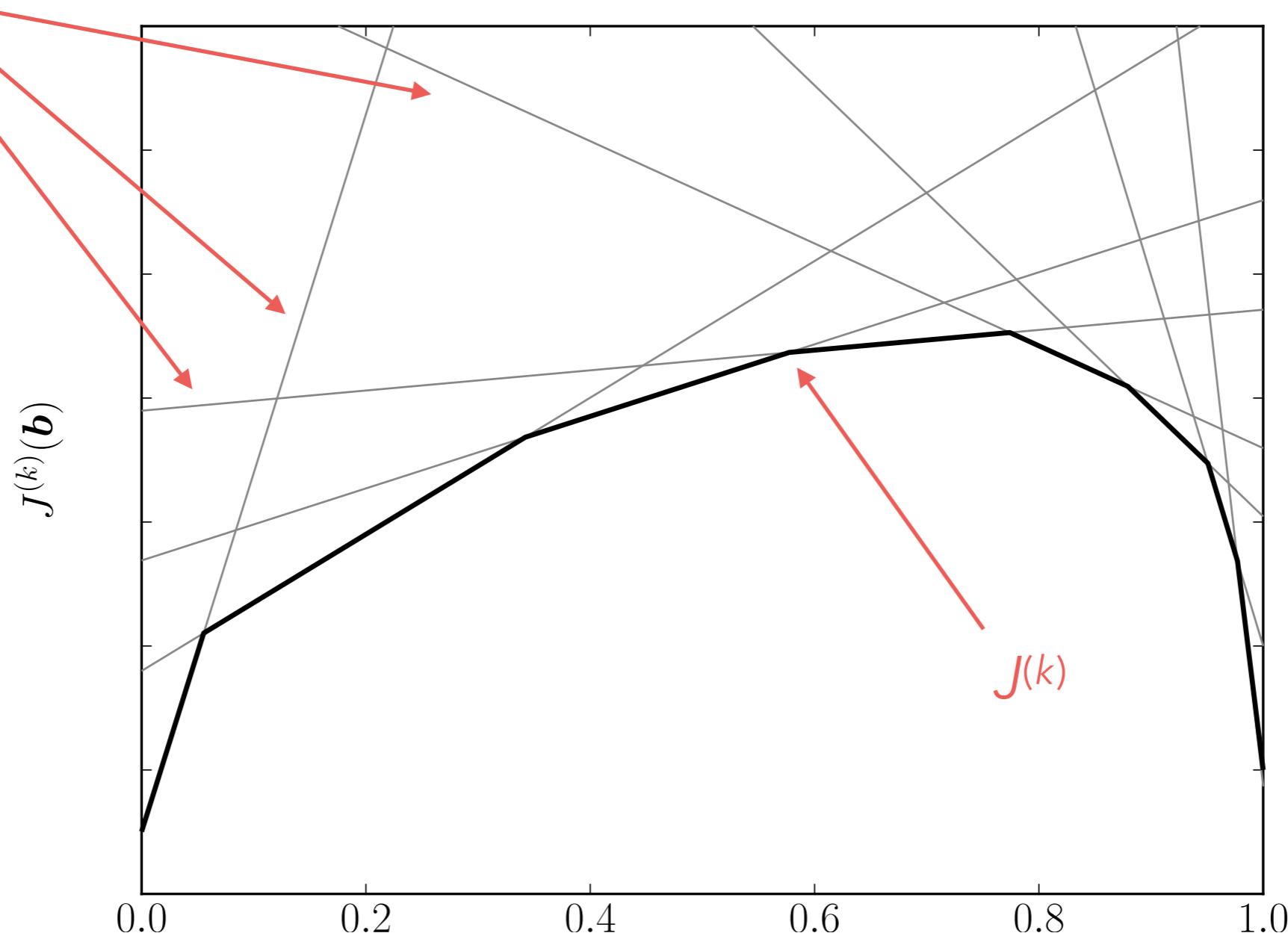
# Representing $J(k)$

*a*-vectors



# Representing $J(k)$

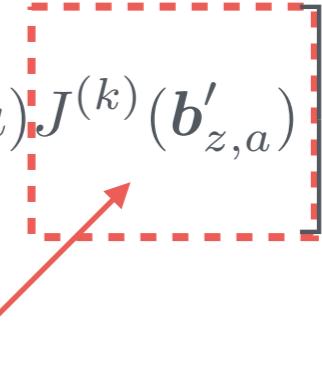
$a$ -vectors



# Value iteration

- How do we compute these a-vectors?

$$J^{(k+1)}(\mathbf{b}) = \min_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) \left[ c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \sum_{y \in \mathcal{X}} \mathbb{P}(y | x, a) \mathbb{O}(z | y, a) J^{(k)}(\mathbf{b}'_{z,a}) \right]$$

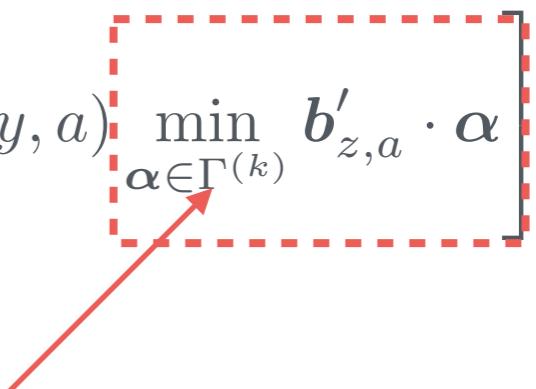


Replace  
representation

# Value iteration

- How do we compute these a-vectors?

$$J^{(k+1)}(\mathbf{b}) = \min_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) \left[ c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \sum_{y \in \mathcal{X}} \mathbb{P}(y | x, a) \mathbb{O}(z | y, a) \min_{\alpha \in \Gamma^{(k)}} \mathbf{b}'_{z,a} \cdot \alpha \right]$$



Minimizing a-vector  
depends on  **$\mathbf{b}$**

# Value iteration

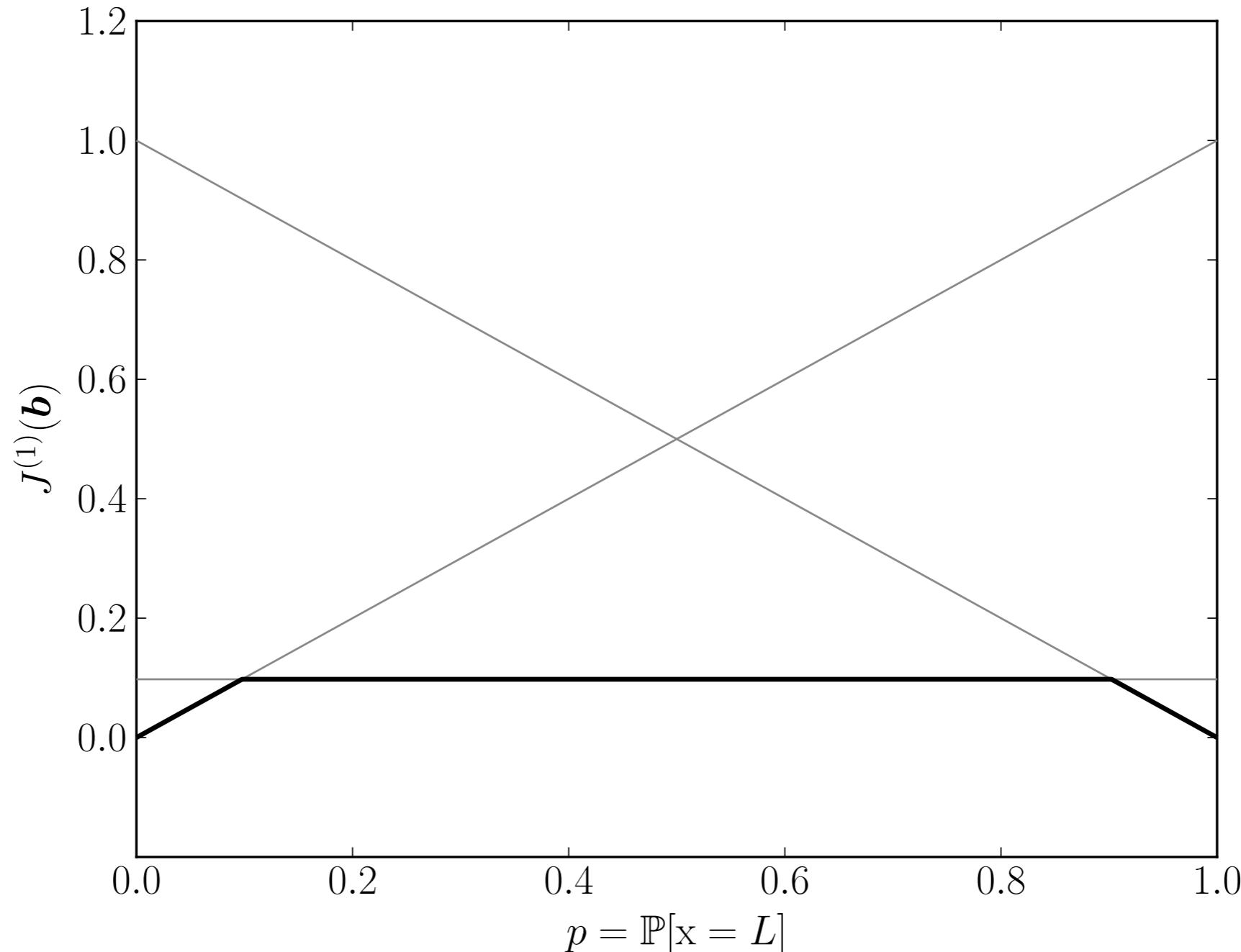
- Compute, at each iteration  $k + 1$ , the set  $\Gamma^{(k+1)}$  from  $\Gamma^{(k)}$ 
  - For each  $\alpha \in \Gamma^{(k)}$ , compute

$$\alpha_{a,z}^{(k)} = \frac{1}{|\mathcal{Z}|} C_{:,a} + \gamma P_a \text{diag}(O_{z,a}) \alpha$$

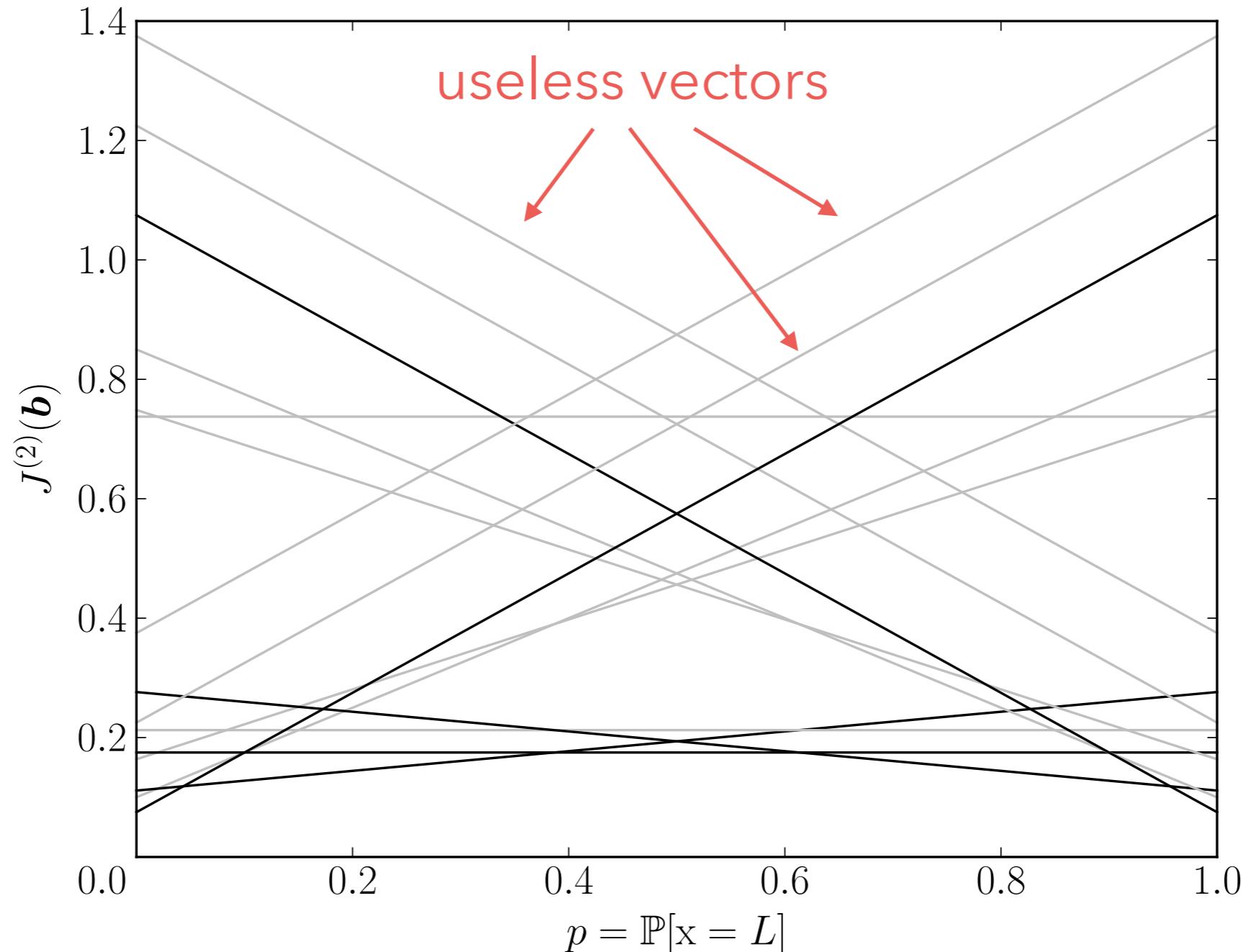
- Compute **all possible combinations** of  $\alpha_{a,z}^{(k)}$ , for each  $z$
- For each combination, let

$$\alpha_a^{(k)} = \sum_{z \in \mathcal{Z}} \alpha_{z,a}^{(k)}$$

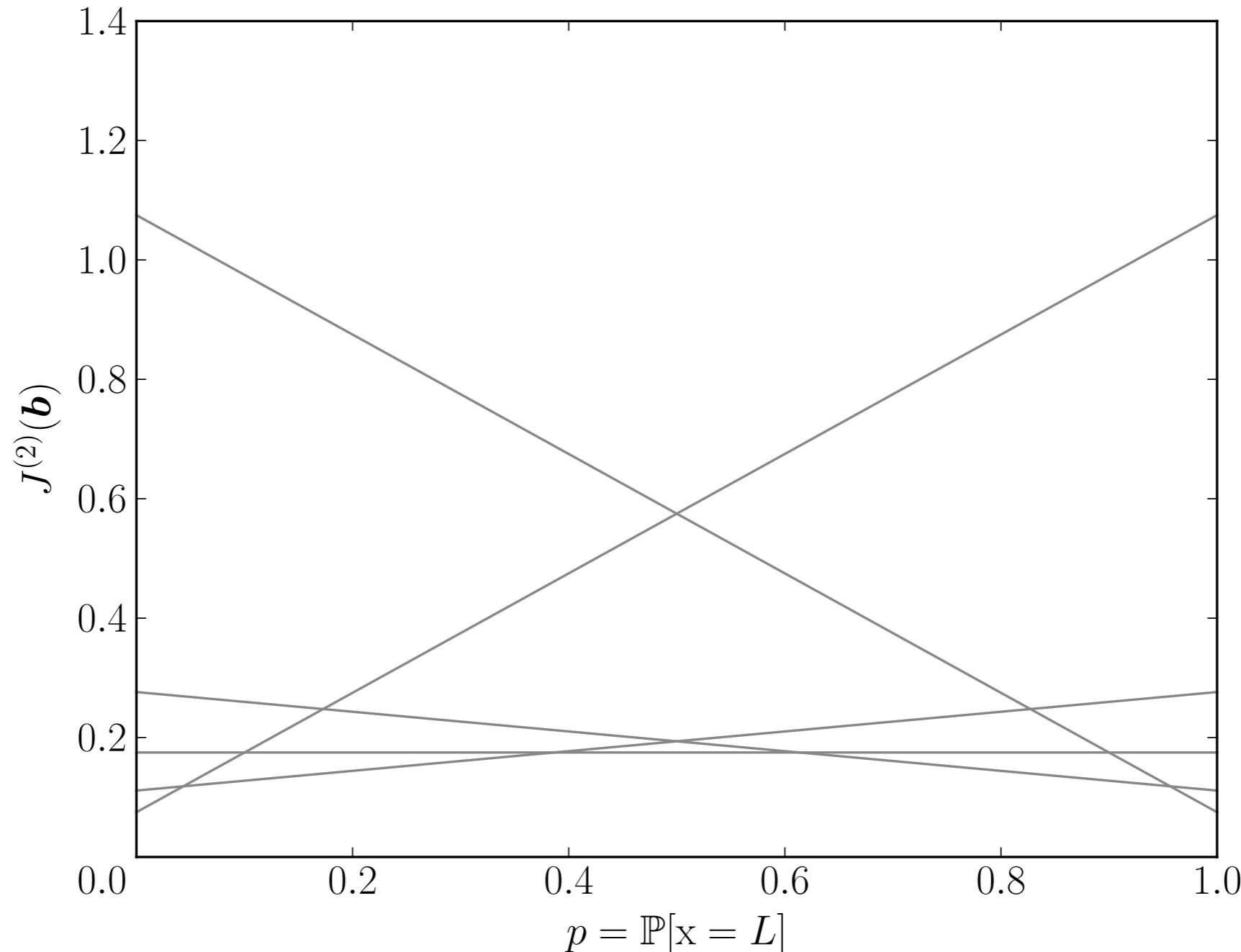
# Value iteration



# Value iteration



# Value iteration



# Value iteration

- Two approaches to build  $\Gamma^{(k+1)}$  from  $\Gamma^{(k)}$ :
  - **Region based methods:** Start with empty  $\Gamma^{(k+1)}$  and only add vectors that are necessary



A vector is necessary if  
it represents  $J$  in a non  
empty belief region  
(witness region)

# Value iteration

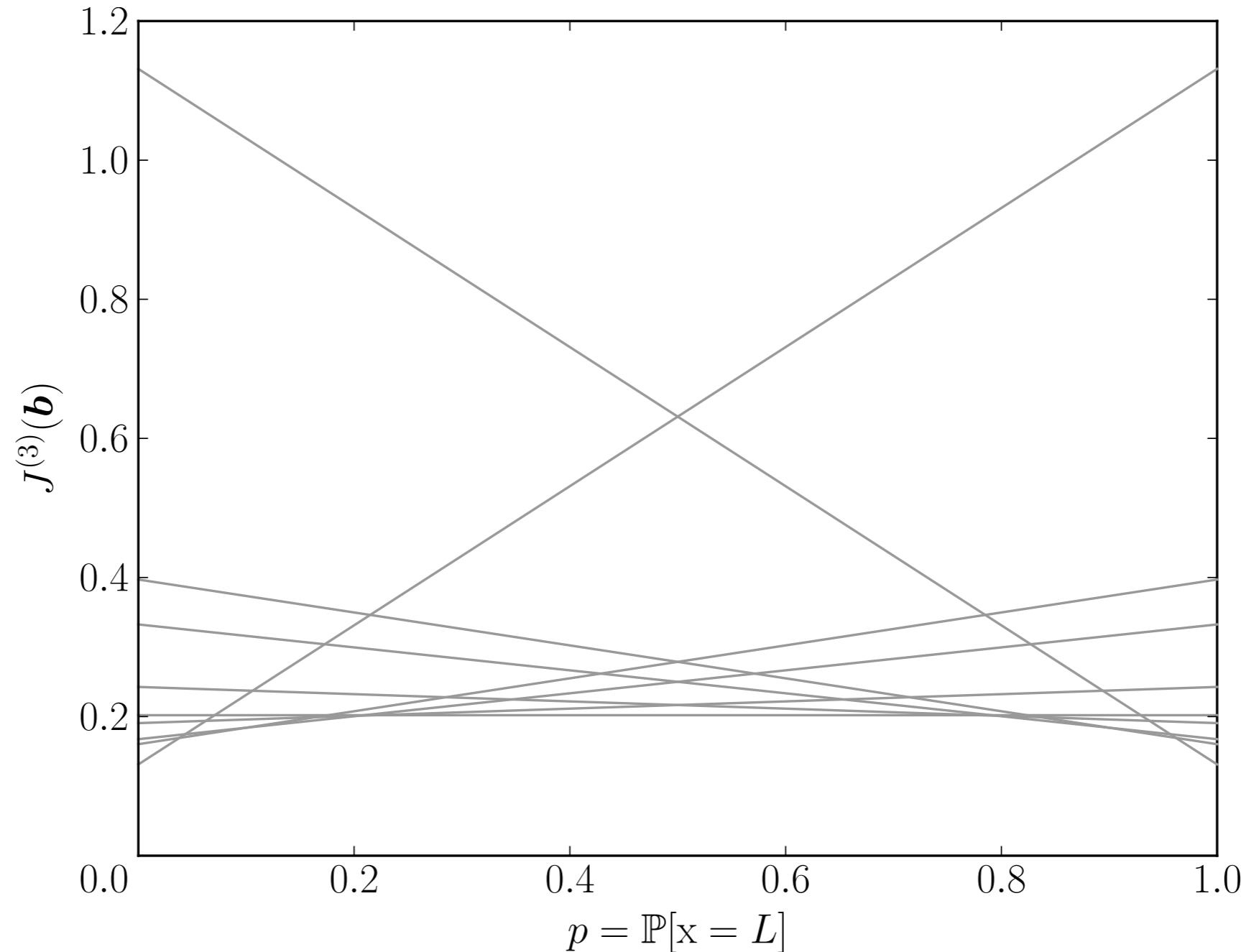
- Two approaches to build  $\Gamma^{(k+1)}$  from  $\Gamma^{(k)}$ :
  - **Region based methods:** Start with empty  $\Gamma^{(k+1)}$  and only add vectors that are necessary

Example: Witness algorithm

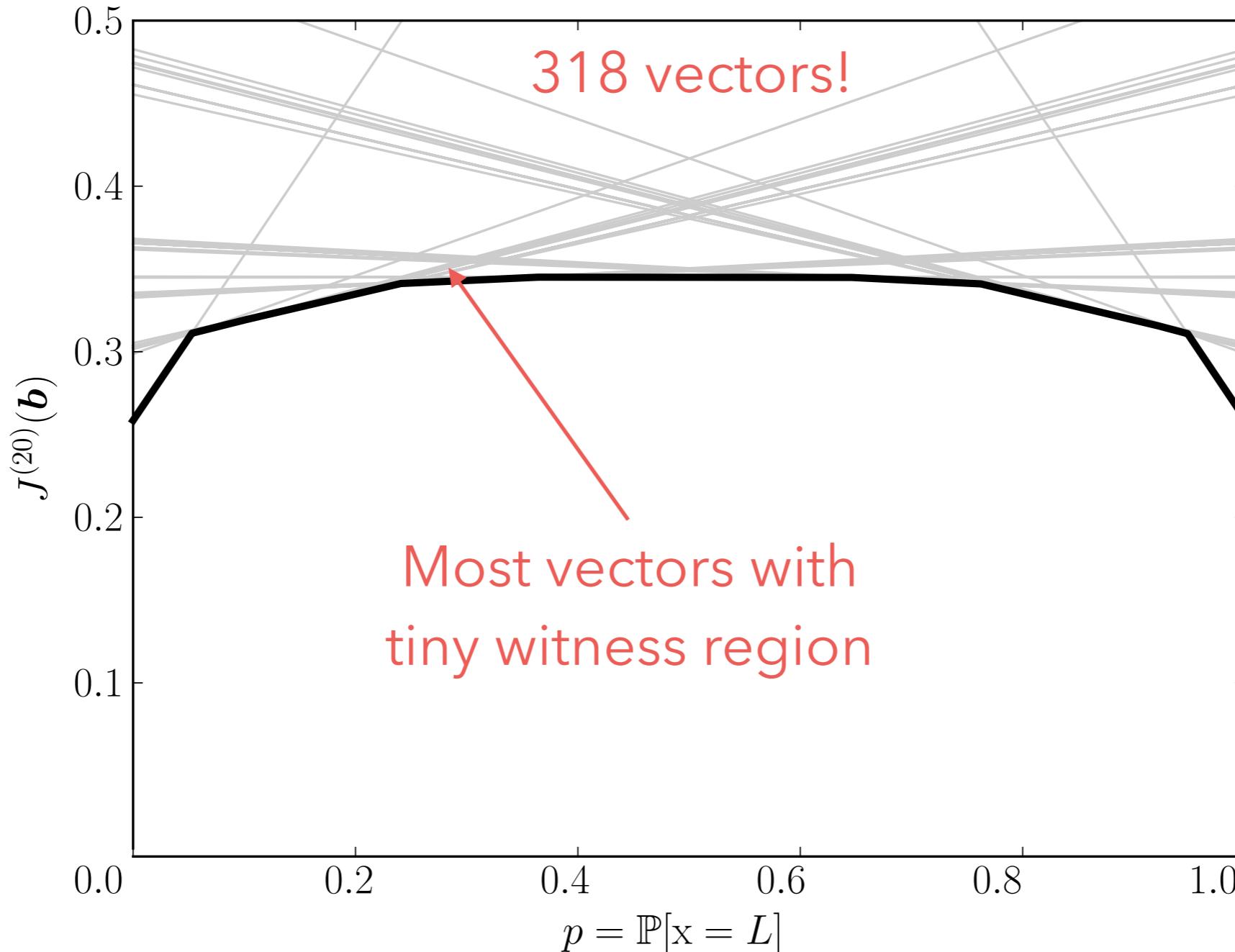
- **Pruning-based methods:** Start with complete  $\Gamma^{(k+1)}$  and remove vectors that are unnecessary

Example: Incremental pruning

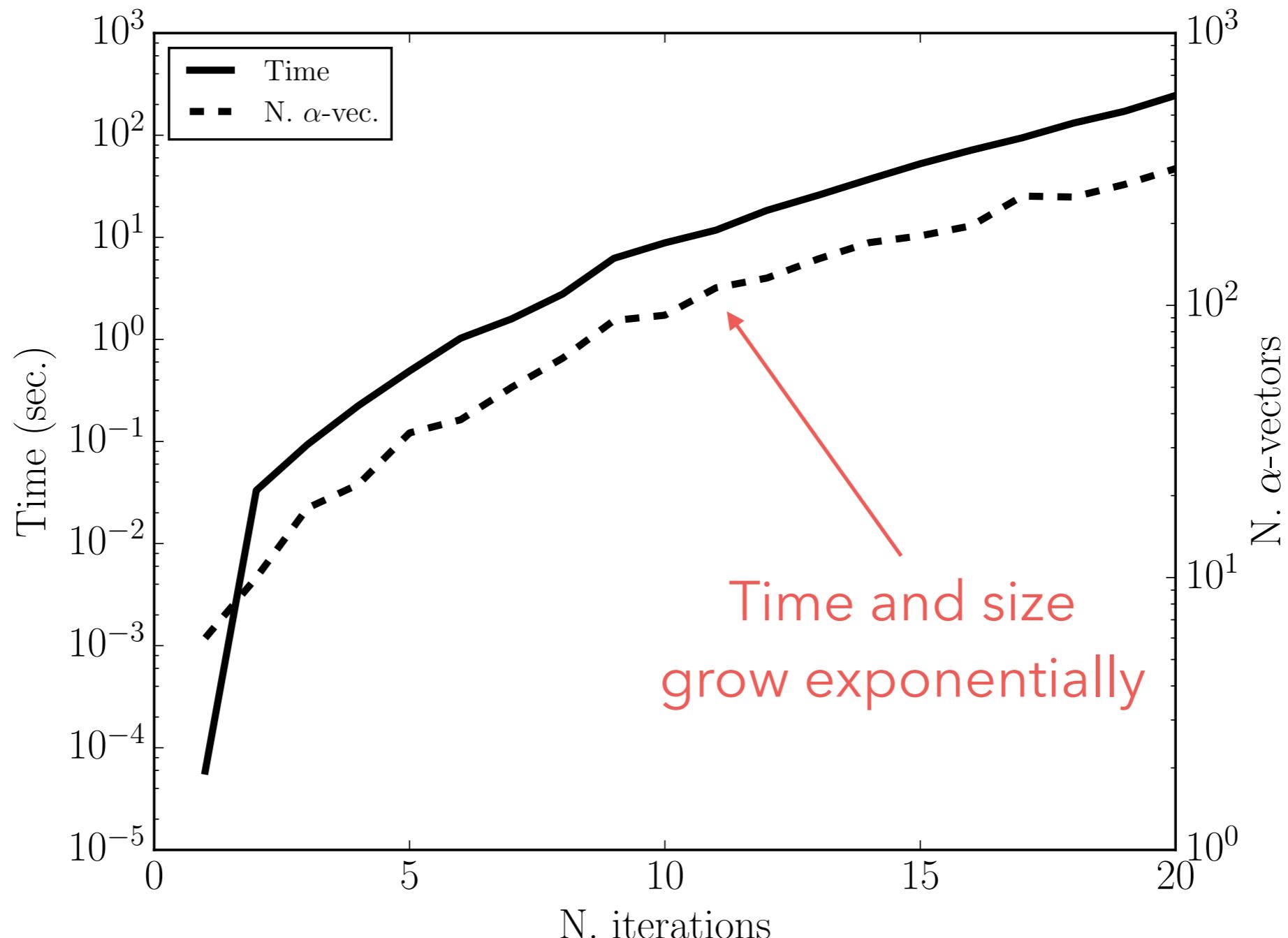
# Value iteration



# Value iteration



# Computation time



# Non-exact solutions

# Idea n. 1 - Use the MDP

- “Under” a POMDP there is an MDP

$$(\mathcal{X}, \mathcal{A}, \mathcal{Z}, \{\mathbf{P}_a\}, \{\mathbf{O}_a\}, c, \gamma)$$

- MDP can be solved efficiently
- Why not use the MDP solution?



MDP heuristics

# MLS heuristic

- Suppose that  $\pi_{MDP}$  is the optimal MDP policy
- At time step  $t$ , we have a belief

$$\mathbf{b} = [\mathbf{b}(x_1) \quad \mathbf{b}(x_2) \quad \dots \quad \mathbf{b}(x_N)]$$



Probability that  
state is  $x_1$

Probability that  
state is  $x_2$

# MLS heuristic

- Suppose that  $\pi_{MDP}$  is the optimal MDP policy
- At time step  $t$ , we have a belief

$$\mathbf{b} = [\mathbf{b}(x_1) \quad \boxed{\mathbf{b}(x_2)} \quad \dots \quad \mathbf{b}(x_N)]$$

Most likely  
state

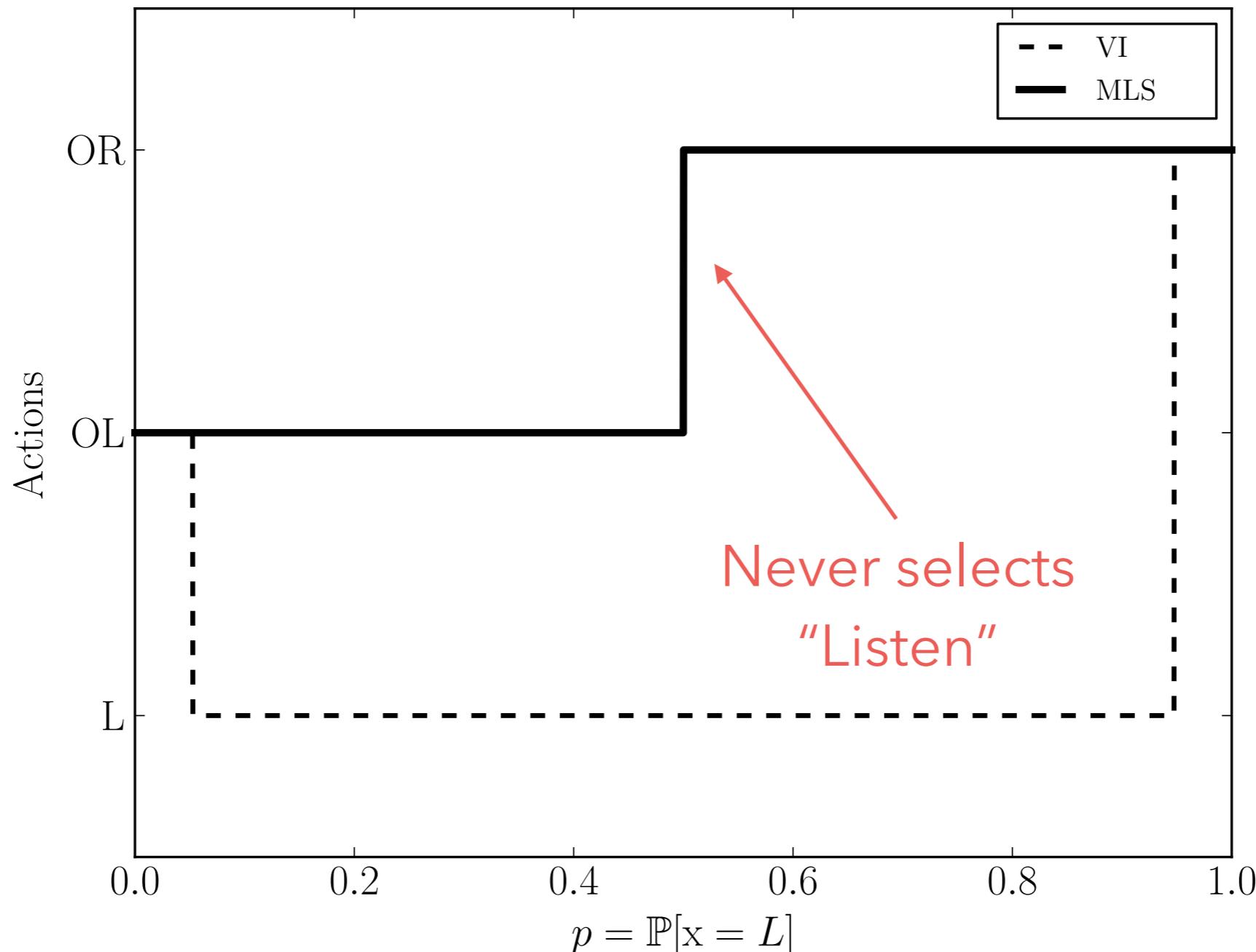
- Select the most likely state (state with largest probability)
- Execute corresponding action –  $\pi_{MDP}(x_2)$

# MLS heuristic

- The MLS (Most Likely State) heuristic is, then

$$\pi_{\text{MLS}}(\mathbf{b}) = \pi_{\text{MDP}}(\operatorname{argmax}_{x \in \mathcal{X}} \mathbf{b}(x))$$

# MLS heuristic



# AV heuristic

- Suppose that  $\pi_{MDP}$  is the optimal MDP policy
- At time step  $t$ , we have a belief

$$\mathbf{b} = [\mathbf{b}(x_1) \quad \mathbf{b}(x_2) \quad \dots \quad \mathbf{b}(x_N)]$$



# AV heuristic

- Suppose that  $\pi_{MDP}$  is the optimal MDP policy

- At time step  $t$ , we have a belief

$$\mathbf{b} = [\mathbf{b}(x_1) \quad \mathbf{b}(x_2) \quad \dots \quad \mathbf{b}(x_N)]$$

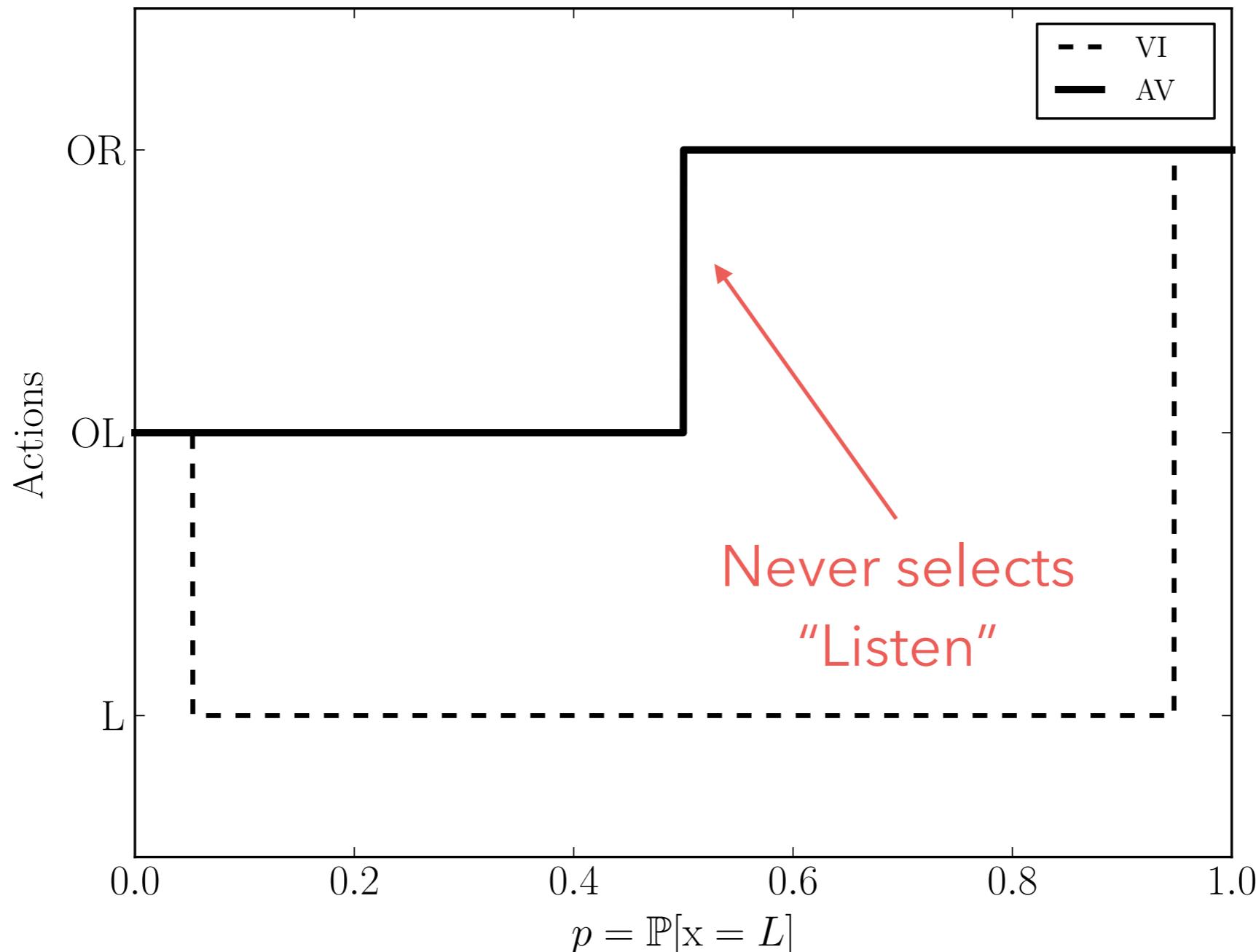
- States vote proportionally to their probability
- Execute most voted action

# AV heuristic

- The AV (Action Voting) heuristic is, then

$$\pi_{\text{AV}}(\mathbf{b}) = \operatorname{argmax}_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) \mathbb{I}(a = \pi_{\text{MDP}}(x))$$

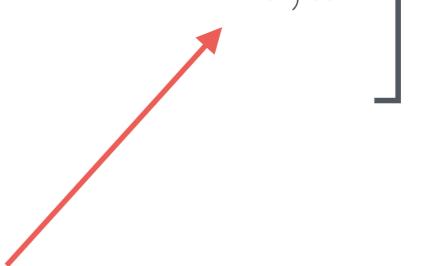
# AV heuristic



# Q-MDP heuristic

- **Optimistic assumption (at time  $t$ ):** partial observability is over at time  $t + 1$
- What is the cost-to-go?

$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) \left[ c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \sum_{y \in \mathcal{X}} P(y | x, a) O(z | y, a) J^*(\mathbf{b}'_{z,a}) \right]$$

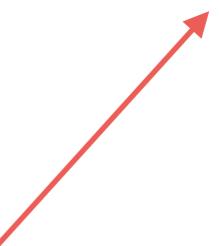


No partial  
observability

# Q-MDP heuristic

- **Optimistic assumption (at time  $t$ ):** partial observability is over at time  $t + 1$
- What is the cost-to-go?

$$J^*(b) = \min_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} b(x) \left[ c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \sum_{y \in \mathcal{X}} P(y | x, a) O(z | y, a) J_{\text{MDP}}(y) \right]$$



No partial  
observability

# Q-MDP heuristic

- **Optimistic assumption (at time  $t$ ):** partial observability is over at time  $t + 1$
- What is the cost-to-go?

$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) \left[ c(x, a) + \gamma \sum_{y \in \mathcal{X}} \mathbb{P}(y | x, a) J_{\text{MDP}}(y) \right]$$

$Q_{\text{MDP}}(x, a)$

# Q-MDP heuristic

- **Optimistic assumption (at time  $t$ ):** partial observability is over at time  $t + 1$
- What is the cost-to-go?

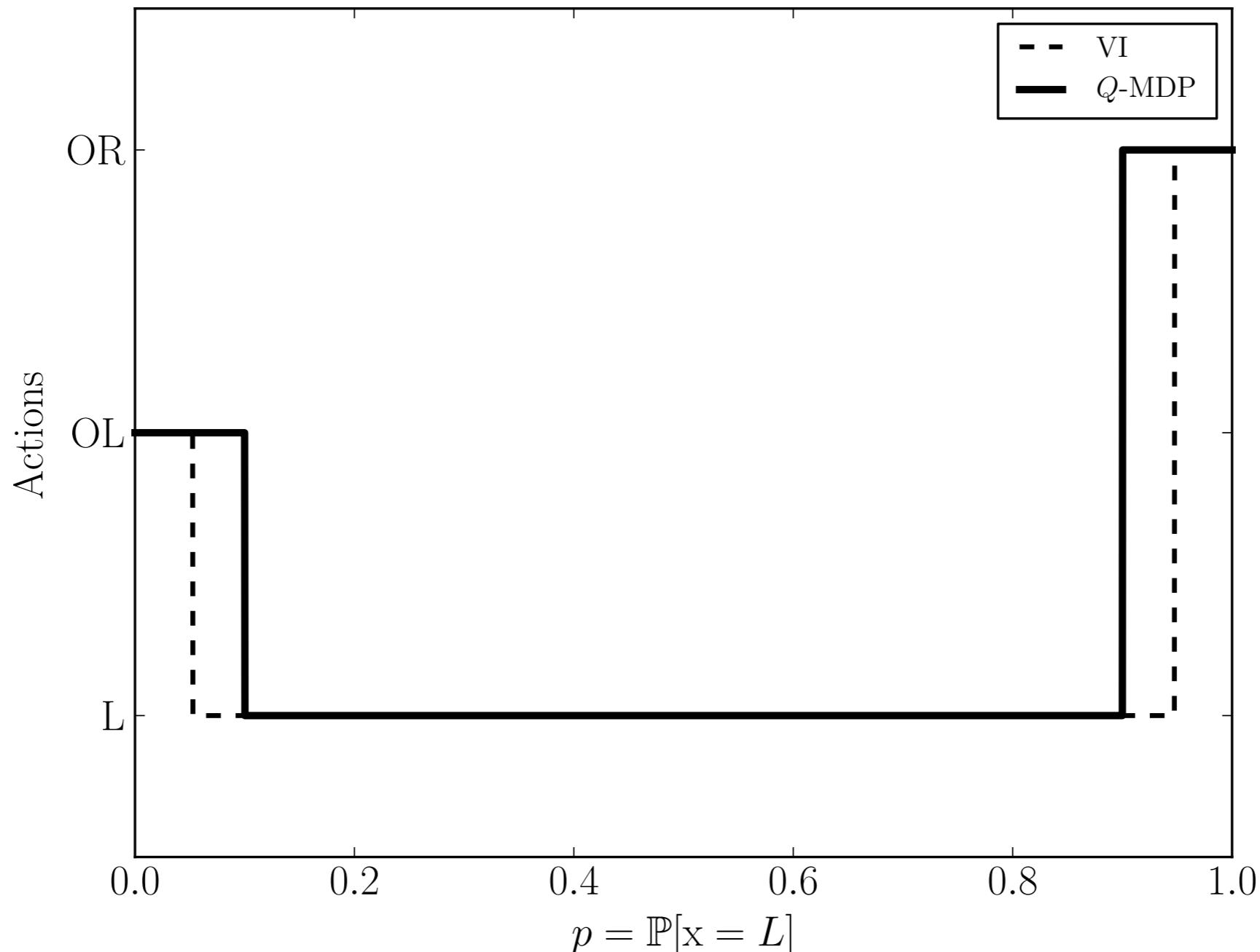
$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) Q_{\text{MDP}}(x, a)$$

# Q-MDP heuristic

- The Q-MDP heuristic is, then

$$\pi_{\text{Q-MDP}}(\mathbf{b}) = \operatorname{argmin}_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) Q_{\text{MDP}}(x, a)$$

# Q-MDP heuristic



# Q-MDP heuristic

- Q-MDP ignores partial observability from the next step on
- What does this mean in terms of cost-to-go approximation?

$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) \left[ c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \sum_{y \in \mathcal{X}} P(y \mid x, a) O(z \mid y, a) \min_{\alpha \in \Gamma^*} \mathbf{b}'_{z,a} \cdot \alpha \right]$$

# Q-MDP heuristic

- Q-MDP ignores partial observability from the next step on
- What does this mean in terms of cost-to-go approximation?

$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \left[ \sum_{x \in \mathcal{X}} \mathbf{b}(x)c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \min_{\alpha \in \Gamma^*} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{X}} \mathbf{b}(x)\mathbb{P}(y | x, a)\mathcal{O}(z | y, a)\alpha_y \right]$$

Optimal

$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \left[ \sum_{x \in \mathcal{X}} \mathbf{b}(x)c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{X}} \mathbf{b}(x)\mathbb{P}(y | x, a)\mathcal{O}(z | y, a) \min_{\alpha \in \Gamma^*} \alpha_y \right]$$

Q-MDP

Ignores partial  
observability



# Q-MDP heuristic

- Q-MDP ignores partial observability from the next step on
- What does this mean in terms of cost-to-go approximation?

$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \left[ \sum_{x \in \mathcal{X}} \mathbf{b}(x)c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \min_{\boldsymbol{\alpha} \in \Gamma^*} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{X}} \mathbf{b}(x)\mathbb{P}(y | x, a)\mathcal{O}(z | y, a)\boldsymbol{\alpha}_y \right]$$

Optimal

$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \left[ \sum_{x \in \mathcal{X}} \mathbf{b}(x)c(x, a) + \gamma \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{X}} \mathbf{b}(x)\mathbb{P}(y | x, a) \min_{\boldsymbol{\alpha} \in \Gamma^*} \boldsymbol{\alpha}_y \right]$$

Q-MDP

# FIB heuristic

- FIB lies somewhere in the middle
- We include partial observability but factor out the belief

$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \left[ \sum_{x \in \mathcal{X}} \mathbf{b}(x) c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \min_{\boldsymbol{\alpha} \in \Gamma^*} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{X}} \mathbf{b}(x) P(y | x, a) O(z | y, a) \boldsymbol{\alpha}_y \right]$$

# FIB heuristic

- FIB lies somewhere in the middle
- We include partial observability but factor out the belief

$$J^*(\mathbf{b}) = \min_{a \in \mathcal{A}} \left[ \sum_{x \in \mathcal{X}} \mathbf{b}(x) c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) \min_{\alpha \in \Gamma^*} \sum_{y \in \mathcal{X}} P(y | x, a) O(z | y, a) \alpha_y \right]$$

# FIB heuristic

- Equivalent to Q-MDP with modified Q-function:

$$Q_{\text{FIB}}(x, a) = c(x, a) + \gamma \sum_{z \in \mathcal{Z}} \min_{a' \in \mathcal{A}} \sum_{y \in \mathcal{X}} P(y \mid x, a) O(z \mid y, a) Q_{\text{FIB}}(y, a')$$

# FIB heuristic

- The FIB heuristic is, then

$$\pi_{\text{FIB}}(\mathbf{b}) = \operatorname{argmax}_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} \mathbf{b}(x) Q_{\text{FIB}}(x, a)$$

# FIB heuristic

