

Romain Hollanders, Raphaël Jungers and Jean-Charles Delvenne

present

**Policy Iteration has exponential complexity  
for solving Markov Decision Processes**

But there is more to the story...

A Cesante Seminar production

**PART 1 : Exponential lower bounds for Policy Iteration**



**PART 2 : Can we obtain good upper bounds ?**

**Romain Hollanders**, Raphaël Jungers and Jean-Charles Delvenne

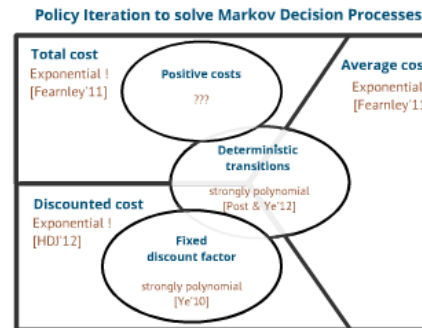
present

**Policy Iteration has exponential complexity  
for solving Markov Decision Processes**

But there is more to the story...

A Cesame Seminar production

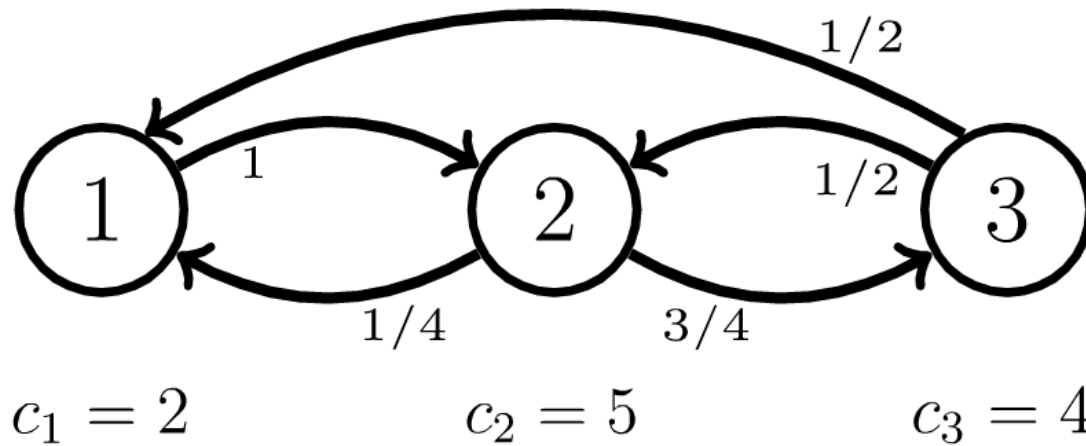
# PART 1 : Exponential lower bounds for Policy Iteration



# PART 2 : Can we obtain good upper bounds ?



# Markov Chains

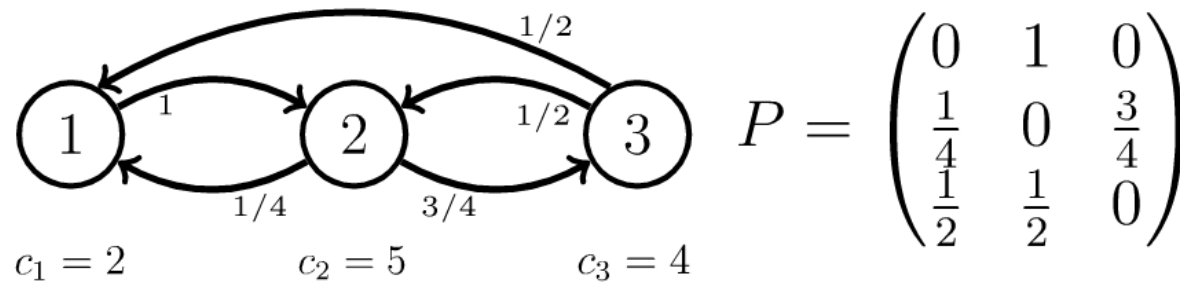


$$P = \begin{pmatrix} 0 & 1 & 0 \\ \frac{1}{4} & 0 & \frac{3}{4} \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix}$$

$$n_i^T = n_0^T P^k$$

How much will I pay

# Markov Chains



$$P = \begin{pmatrix} 0 & 1 & 0 \\ \frac{1}{4} & 0 & \frac{3}{4} \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix}$$

$$p_k^T = p_0^T P^k$$

$$p_0^T = 1 \quad 0 \quad 0$$

$$p_1^T = 0 \quad 1 \quad 0$$

$$p_2^T = \frac{3}{4} \quad 0 \quad \frac{1}{4}$$

$$p_3^T = \frac{1}{8} \quad \frac{7}{8} \quad 0$$

$$p_4^T = \frac{21}{32} \quad \frac{4}{32} \quad \frac{7}{32}$$

⋮

How much will I pay if I start from state 1?

- Total cost

$$x(1) = \sum_{k=0}^H p_k^T c$$

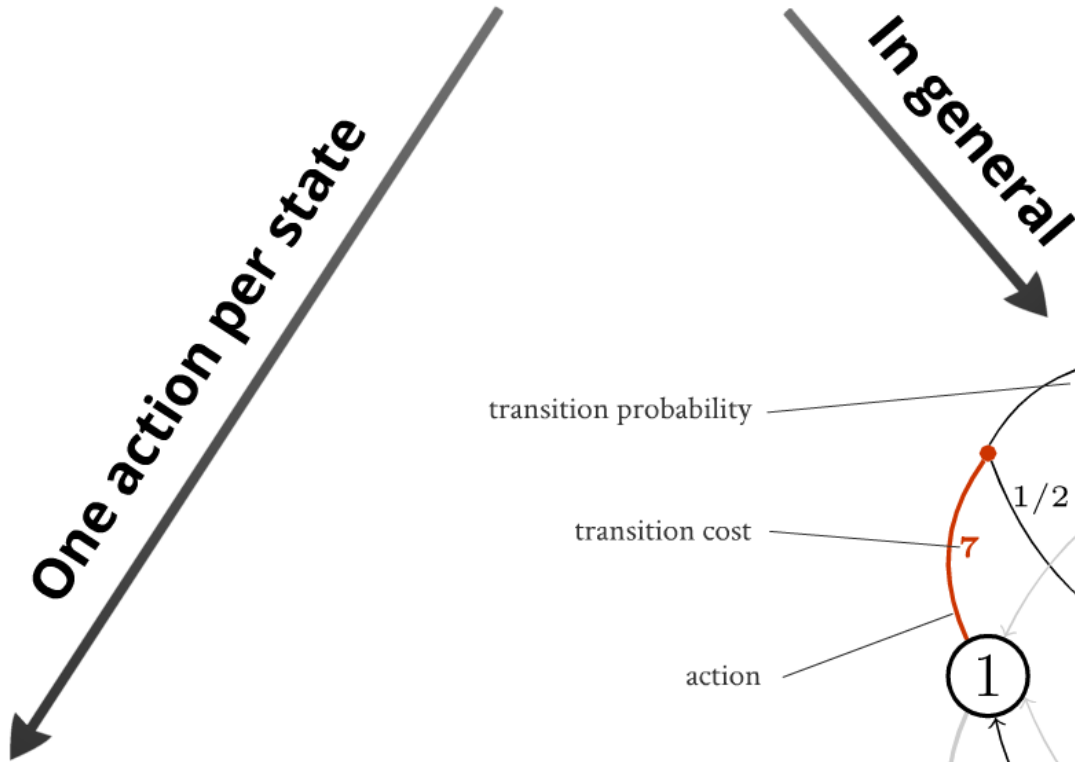
- Average cost

$$x(1) = \lim_{H \rightarrow \infty} \frac{1}{H} \cdot \sum_{k=0}^H p_k^T c$$

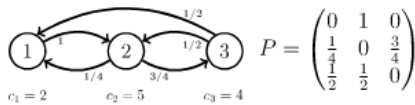
- Discounted cost

$$x(1) = \sum_{k=0}^{\infty} \gamma^k p_k^T c$$

# Markov Decision Processes



## Markov Chains



$$p_k^T = p_0^T P^k$$

$$p_0^T = 1 \ 0 \ 0$$

$$p_1^T = 0 \ 1 \ 0$$

$$p_2^T = \frac{3}{4} \ 0 \ \frac{1}{4}$$

$$p_3^T = \frac{1}{8} \ \frac{7}{8} \ 0$$

$$p_4^T = \frac{21}{32} \ \frac{4}{32} \ \frac{7}{32}$$

$$\vdots$$

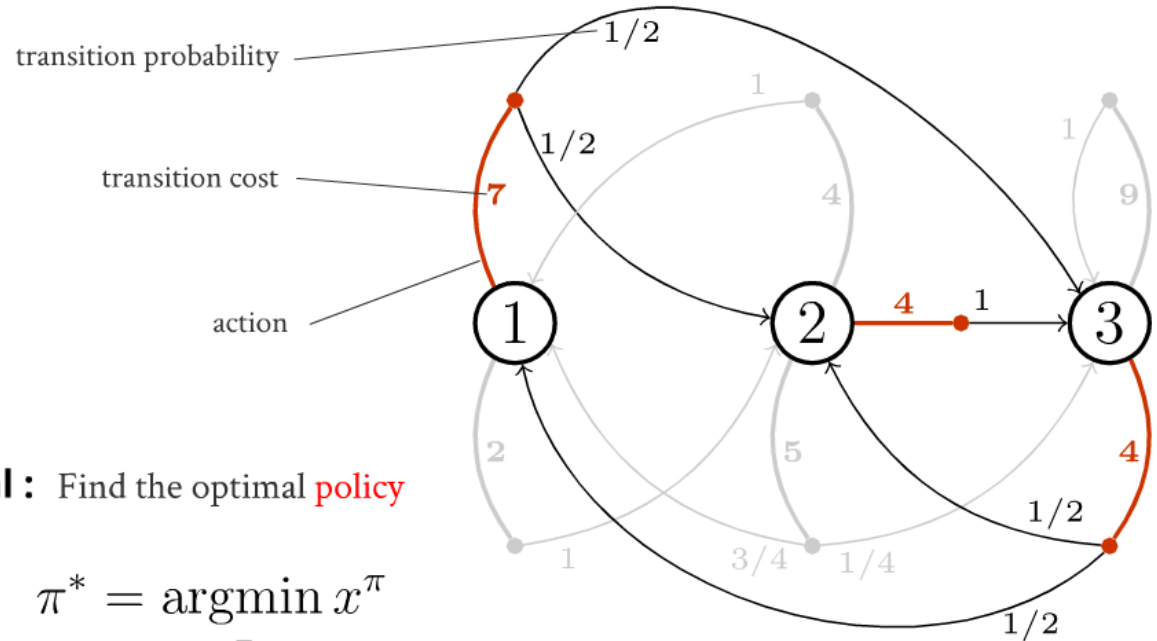
How much will I pay if I start from state 1?

- Total cost  

$$x(1) = \sum_{k=0}^H p_k^T c$$
- Average cost  

$$x(1) = \lim_{H \rightarrow \infty} \frac{1}{H} \cdot \sum_{k=0}^H p_k^T c$$
- Discounted cost  

$$x(1) = \sum_{k=0}^{\infty} \gamma^k p_k^T c$$

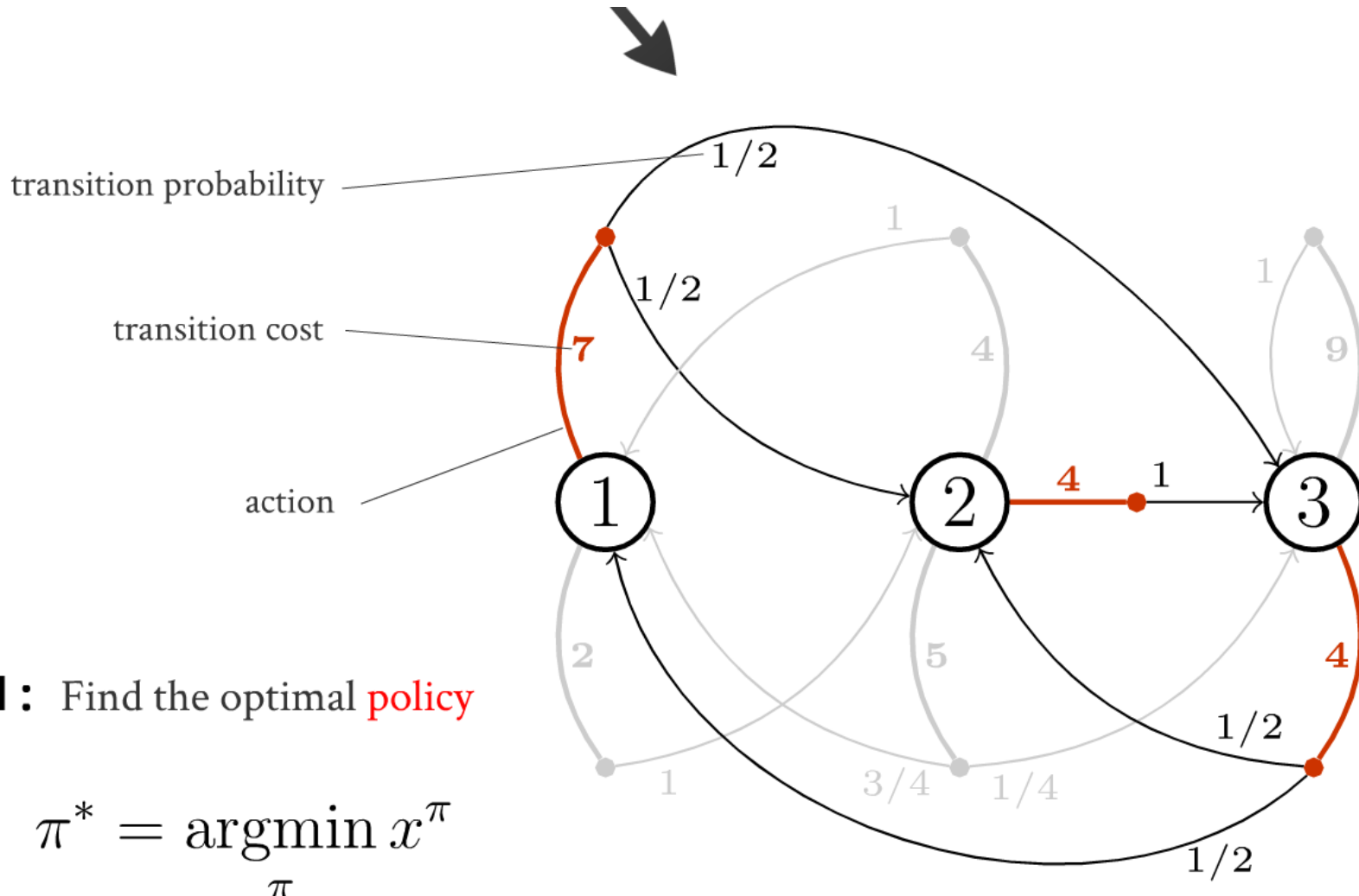


**Goal :** Find the optimal **policy**

$$\pi^* = \operatorname{argmin}_{\pi} x^{\pi}$$

The answer depends on the chosen objective function :

- Total cost
- Average cost
- Discounted cost



**Goal :** Find the optimal **policy**

$$\pi^* = \operatorname{argmin}_{\pi} x^{\pi}$$

The answer depends on the chosen objective function :

- Total cost
- Average cost
- Discounted cost

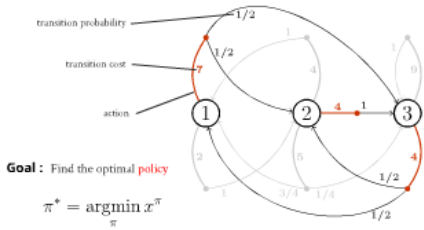
# Markov Decision Processes

One action per state

In general

### Markov Chains

$x^k = P^k x^0$   
 How much will per  
 if I start from state 1?  
 • Total cost  
 $J^k = \sum_{t=0}^{\infty} c^t$   
 • Average cost  
 $J^k = \lim_{k \rightarrow \infty} \frac{1}{k} \sum_{t=0}^{k-1} c^t$   
 • Discounted cost  
 $J^k = \sum_{t=0}^{\infty} \gamma^t c^t$



Goal: Find the optimal policy

$$\pi^* = \operatorname{argmin}_{\pi} x^{\pi}$$

The answer depends on the chosen objective function:

- Total cost
- Average cost
- Discounted cost

Resolution

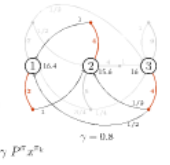
# Linear Programming

Value Iteration

Policy Iteration

```

0. Choose an initial policy  $\pi_0$ 
while  $\pi_k \neq \pi_{k-1}$ 
1. Evaluate  $\pi_k$  (Bellman)
 $x^{\pi_k} = c^{\pi_k} + \gamma P^{\pi_k} x^{\pi_k}$ 
2. Improve  $\pi_k$ 
 $\pi_{k+1} = \operatorname{argmin}_{\pi} c^{\pi} + \gamma P^{\pi} x^{\pi_k}$ 
 $k \leftarrow k + 1$ 
end
    
```





# Linear Programming

## Value Iteration

## Policy Iteration

0. Choose an initial policy  $\pi_0$

while  $\pi_k \neq \pi_{k-1}$

1. Evaluate  $\pi_k$  (Bellman)

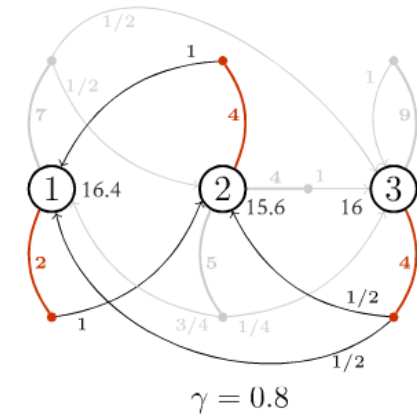
$$x^{\pi_k} = c^{\pi_k} + \gamma P^{\pi_k} x^{\pi_k}$$

2. Improve  $\pi_k$

$$\pi_{k+1} = \operatorname{argmin}_{\pi} c^{\pi} + \gamma P^{\pi} x^{\pi_k}$$

$k \leftarrow k + 1$

end



# Policy Iteration

0. Choose an initial policy  $\pi_0$

while  $\pi_k \neq \pi_{k-1}$

1. Evaluate  $\pi_k$  (Bellman)

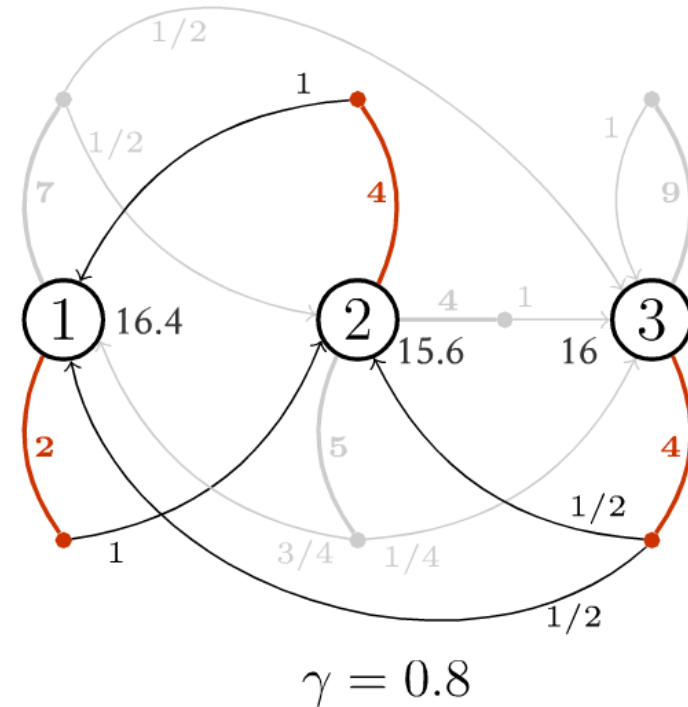
$$x^{\pi_k} = c^{\pi_k} + \gamma P^{\pi_k} x^{\pi_k}$$

2. Improve  $\pi_k$

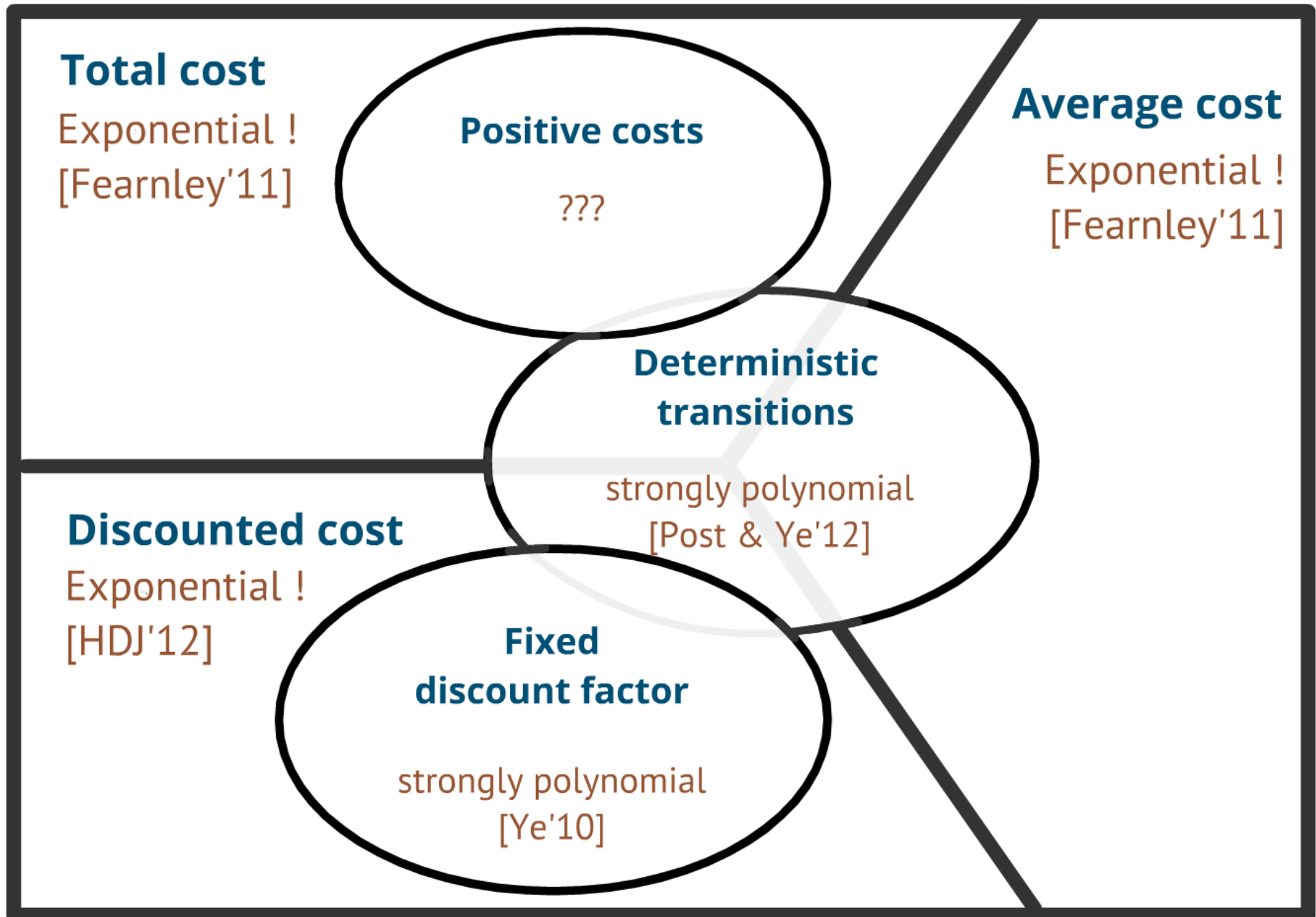
$$\pi_{k+1} = \operatorname{argmin}_{\pi} c^{\pi} + \gamma P^{\pi} x^{\pi_k}$$

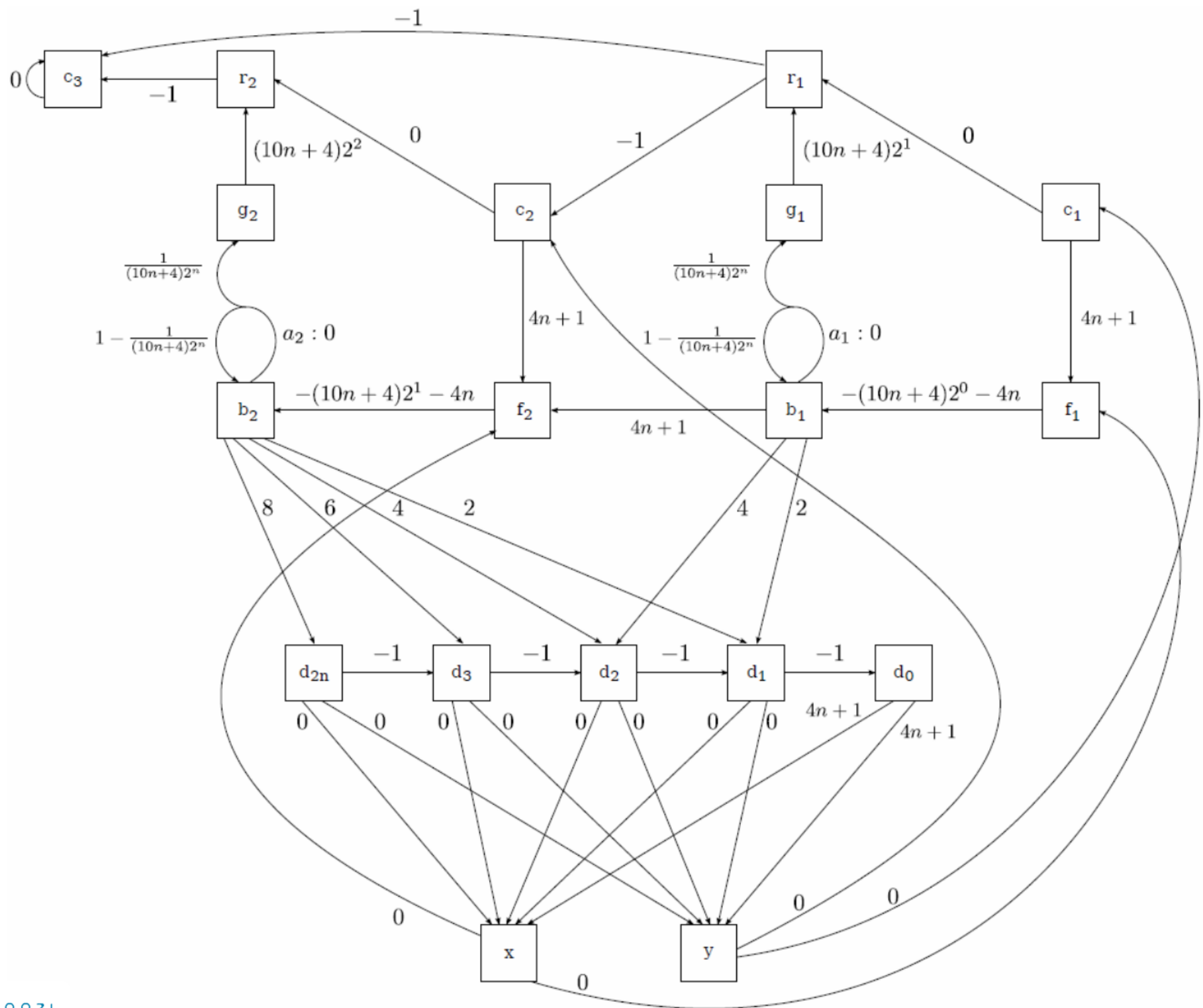
$k \leftarrow k + 1$

end

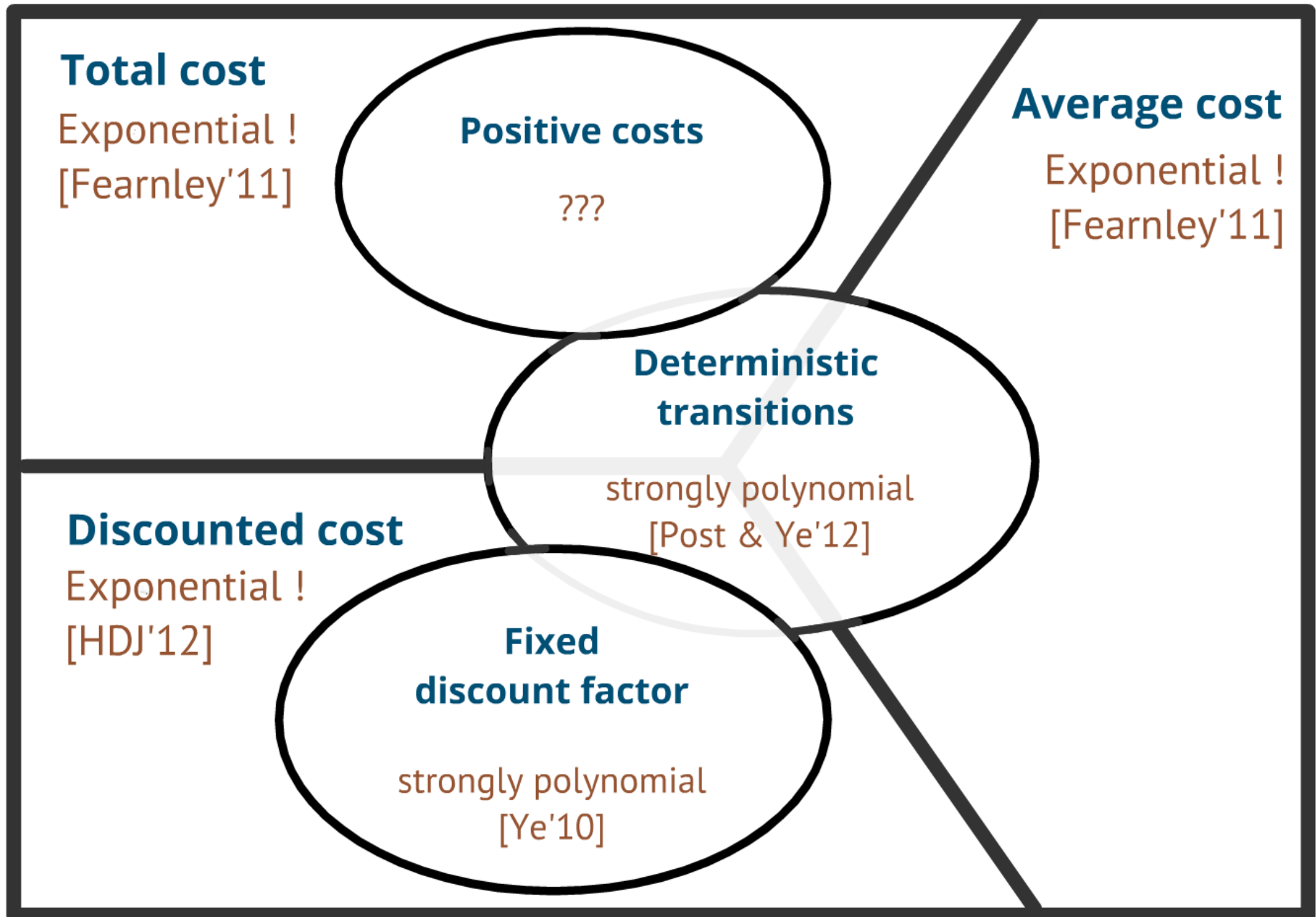


# Policy Iteration to solve Markov Decision Processes



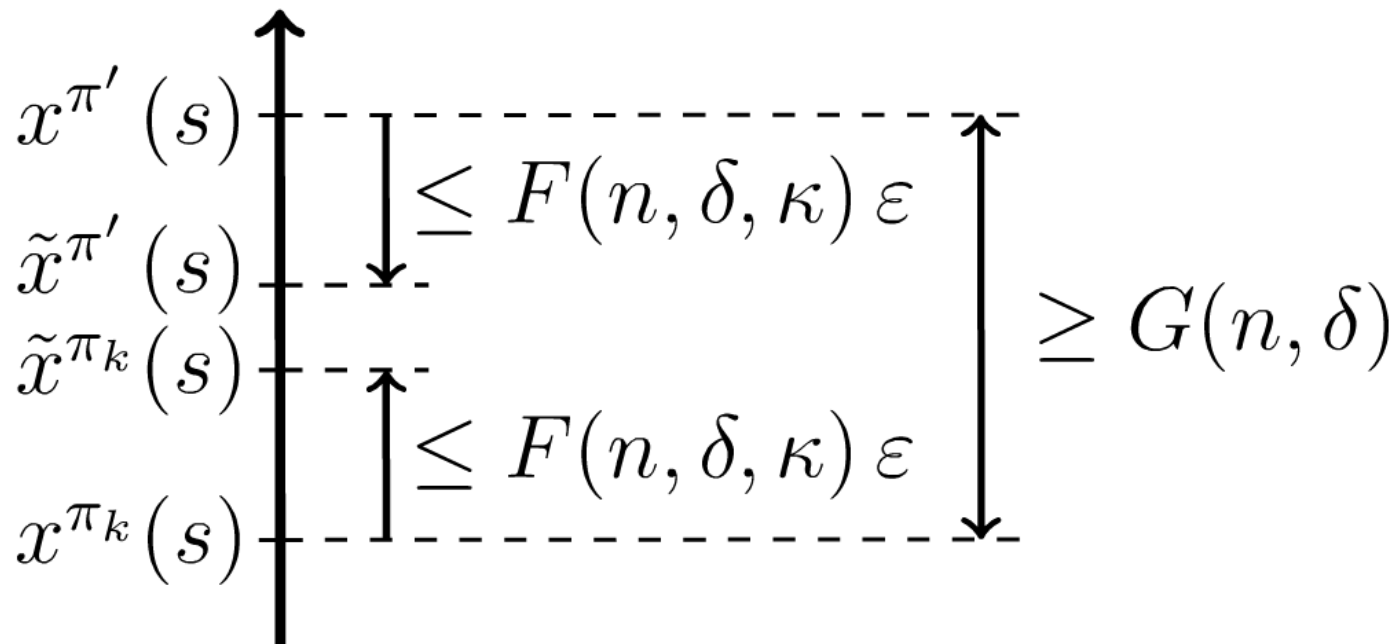
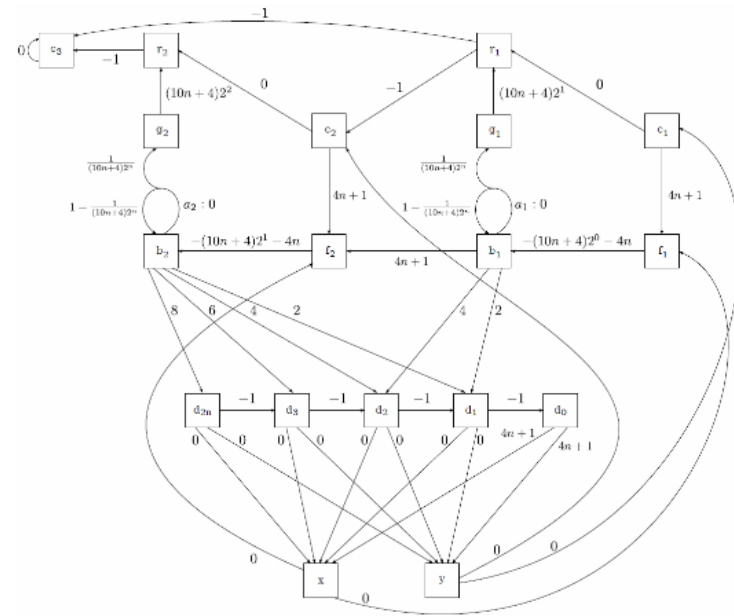


# Policy Iteration to solve Markov Decision Processes



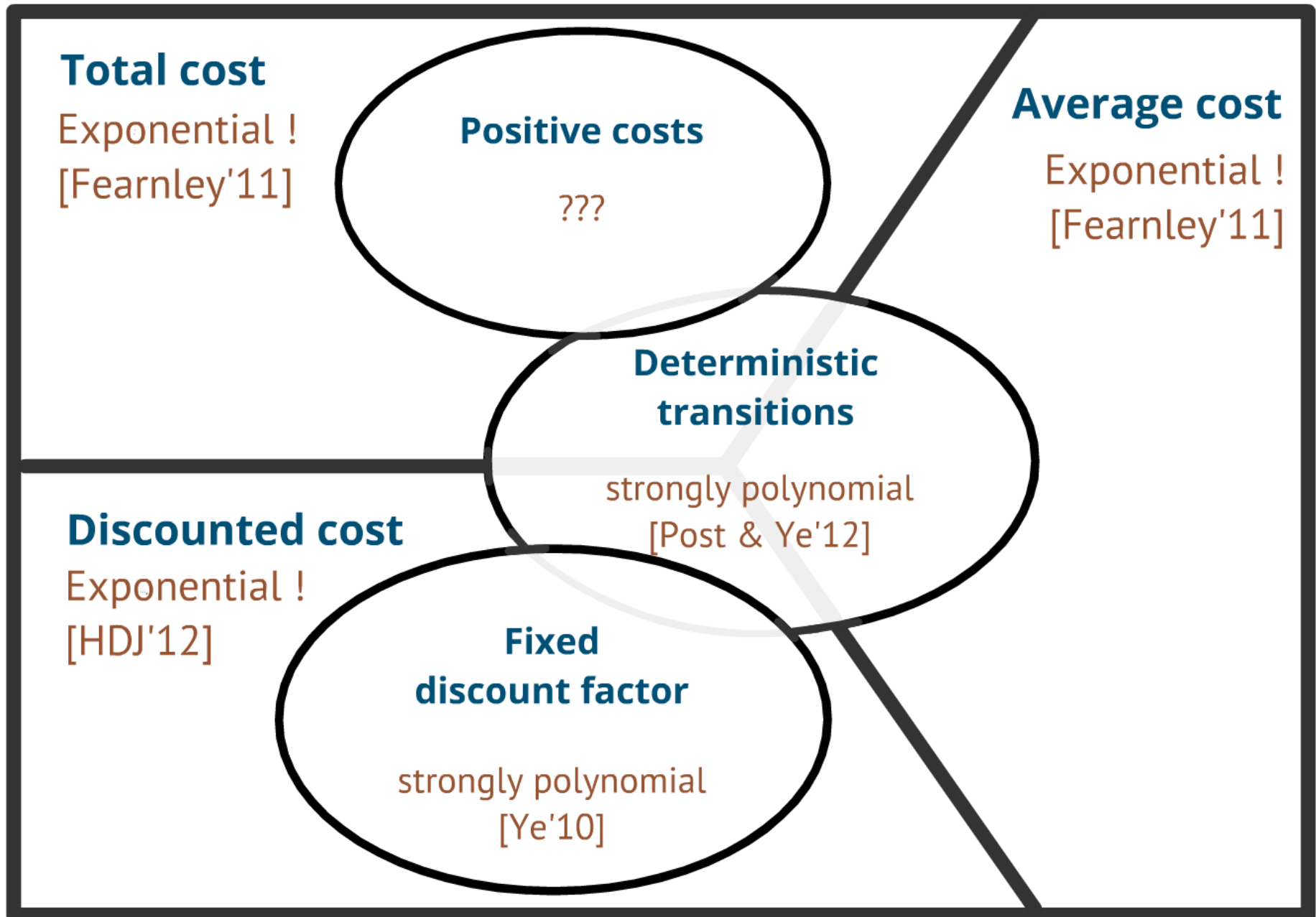
We add discount  $\gamma = 1 - \varepsilon$

➔ How much perturbation?

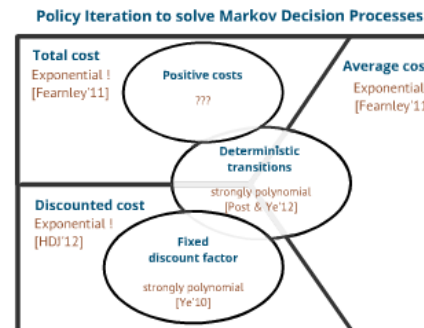


➔ OK for some  $\varepsilon \sim \frac{1}{2q(n, \delta, \kappa)}$

# Policy Iteration to solve Markov Decision Processes



# PART 1 : Exponential lower bounds for Policy Iteration



# PART 2 : Can we obtain good upper bounds ?



# PART 1 : Exponential lower bounds for Policy Iteration



# PART 2 : Can we obtain good upper bounds?

Theorem of Fiacchi and Santopietro '16

Approximation	$\epsilon^2 \cdot n^2 \cdot \log(1/\epsilon)$
Complexity	$\epsilon^{-2} \cdot n^2 \cdot \log(1/\epsilon)$
Best known upper	$\epsilon^{-2} \cdot n^2 \cdot \log(1/\epsilon)$

Linear Programming

**Simplex**

$\Omega(\text{Exponential})$   
[Friedmann '11]  
[Friedmann, Hansen & Zwick '17]

$\Pi =$

0	0	0
1	1	1
1	0	0
1	1	0
0	1	1

Markov Decision Processes

**Policy Iteration**

$\Omega(\text{Exponential})$   
[Fearnley '10]  
 $O(2^n/n)$   
[Mansour & Singh '99]

$\Pi =$

0	0	0
1	1	1
1	0	0
1	1	0
0	1	0

**Acyclic Unique Sink Orientations**

Every face has a unique sink  
No cycle

**The goal:**  
Bound the length of any "PI-path"  
in any Acyclic Unique Sink Orientation

**Order-Regular matrices**

$\Pi \in \{0,1\}^{m \times n}$  is Order-Regular  
iff  $\forall 0 < i < j < m, \exists 0 < k \leq n, s.t.$   
 $\Pi_{i,k} \neq \Pi_{i+1,k} = \Pi_{i,k} = \Pi_{j+1,k}$

**The goal:**  
Bound the number of rows  
of any Order-Regular matrix

**Best known upper bound:**  
 $O(2^n/n)$  [Mansour & Singh '99]

**Experimentally:**

<b>n</b>	1	2	3	4	5	6	7
<b>m*</b>	2	3	5	8	13	21	34

**Conjecture** [Hansen & Zwick '12]

$m^* = F_{n+2}$   
 $= O(\phi^n)$   
golden ratio

# Summary on the complexity of Policy Iteration

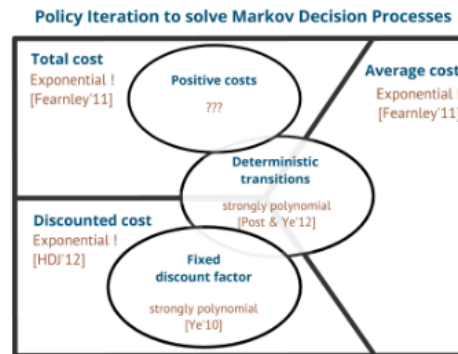
**Lower bounds**

- In general: **exponential** in the worst case
- With some particular structures: **strongly polynomial**
- What about the **average complexity**?

**Upper bounds**

- Best known bound:  $O(2^n/n)$
- Our hope:  $O(1.618^n)$

# PART 1 : Exponential lower bounds for Policy Iteration



# PART 2 : Can we obtain good upper bounds?

The state of the art and some inspiring ideas

Upper bounds on  $m^*$        $m^* \leq O(2^n/n) \leq O(2^n)$

Conjecture on  $m^*$        $m^* \sim \Theta(\phi^n) = O(1.618^n)$

Lower bounds on  $m^*$        $m^* \geq \Omega(\sqrt{2}^n) = \Omega(1.4142^n)$

objective  
function  
→

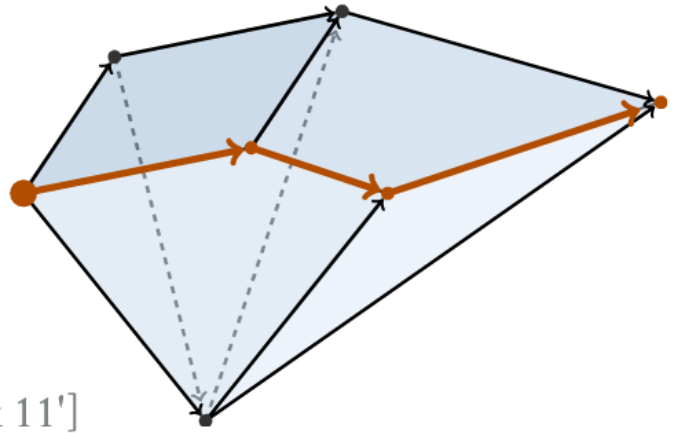
**Linear  
Programming**

**Markov  
Decision  
Processes**



### Simplex

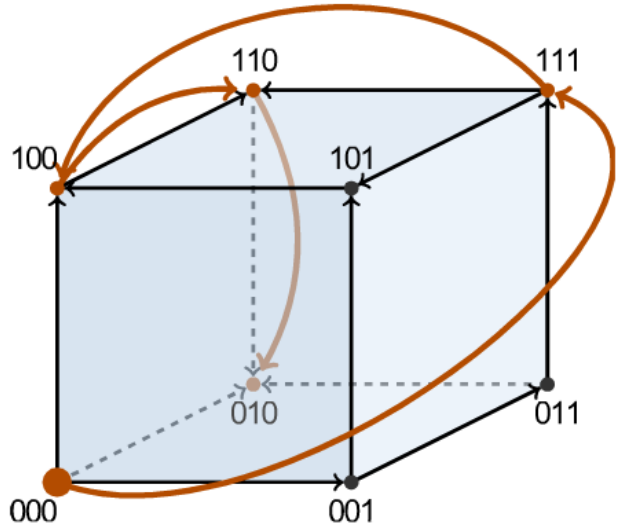
$\Omega(\text{Exponential})$   
[Friedmann '11]  
[Friedmann, Hansen & Zwick 11']

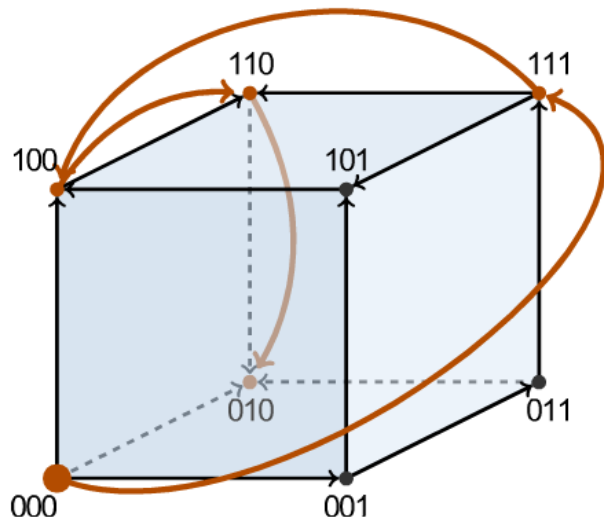


### Policy Iteration

$\Omega(\text{Exponential})$   
[Fearnley '10]

$O(2^n/n)$   
[Mansour & Singh '99]





0	1	1
---	---	---


 $\Pi =$ 

0	0	0
1	1	1
1	0	0
1	1	0
0	1	0



## Acyclic Unique Sink Orientations

Every face has a unique sink

No cycle

### The goal:

Bound the length of any "PI-path"  
in any Acyclic Unique Sink Orientation

## Order-Regular matrices

$\Pi \in \{0, 1\}^{m \times n}$  is Order-Regular

iff  $\forall 0 < i < j < m, \exists 0 < k \leq n$  s.t.:

$$\Pi_{i,k} \neq \Pi_{i+1,k} = \Pi_{j,k} = \Pi_{j+1,k}$$

### The goal:

Bound the number of rows  
of any Order-Regular matrix

Best

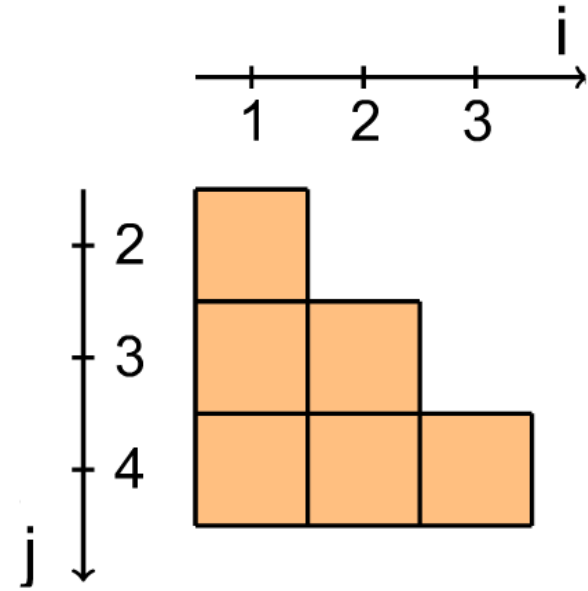
O

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



$\Pi =$

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



## Order-Regular matrices

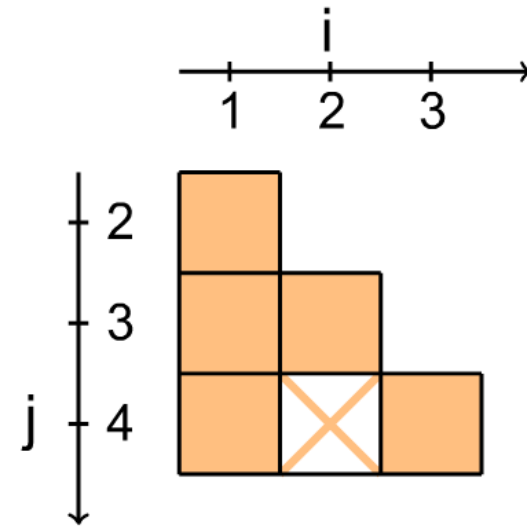
$\Pi \in \{0, 1\}^{m \times n}$  is Order-Regular

iff  $\forall 0 < i < j < m, \exists 0 < k \leq n$  s.t.:

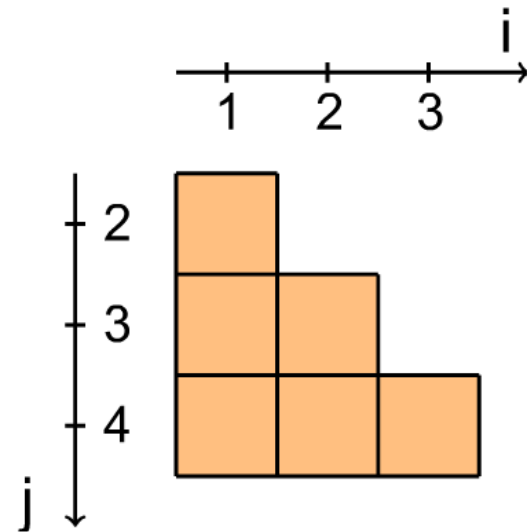
$$\Pi_{i,k} \neq \Pi_{i+1,k} = \Pi_{j,k} = \Pi_{j+1,k}$$

**Best known upper bound:**

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



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



$$\Pi_{i,k} \neq \Pi_{i+1,k} = \Pi_{j,k} = \Pi_{j+1,k}$$

## The goal:

Bound the number of rows  
of any Order-Regular matrix

## Best known upper bound:

$$O(2^n/n) \text{ [Mansour \& Singh '99]}$$

## Experimentally:

<b>n</b>	1	2	3	4	5	6	7
<b>m*</b>	2	3	5	8	13	21	34

## Conjecture [Hansen & Zwick '12]

$$\begin{aligned} m^* &= F_{n+2} \\ &= O(\phi^n) \end{aligned}$$

golden ratio

# The state of the art and some inspiring ideas

**Upper bounds on  $m^*$**

$$m^* \leq O(2^n/n) \leq O(2^n)$$

**Conjecture on  $m^*$**

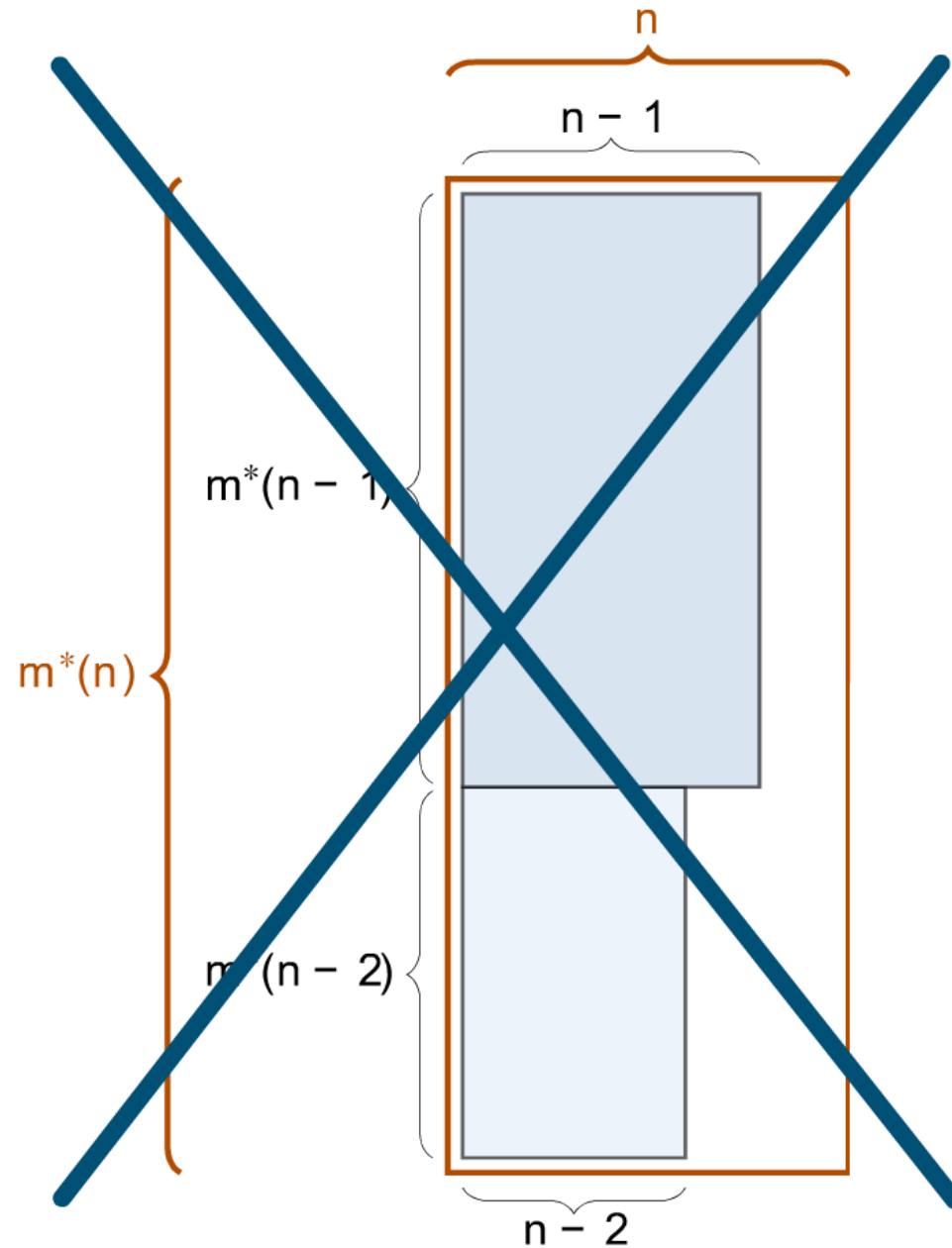
$$m^* \sim O(\phi^n) = O(1.618^n)$$

**Lower bounds on  $m^*$**

$$m^* \geq \Omega(\sqrt{2}^n) = \Omega(1.4142^n)$$



# A natural idea



**fails...**

# The state of the art and some inspiring ideas

**Upper bounds on  $m^*$**

$$m^* \leq O(2^n/n) \leq O(2^n)$$

**Conjecture on  $m^*$**

$$m^* \sim O(\phi^n) = O(1.618^n)$$

**Lower bounds on  $m^*$**

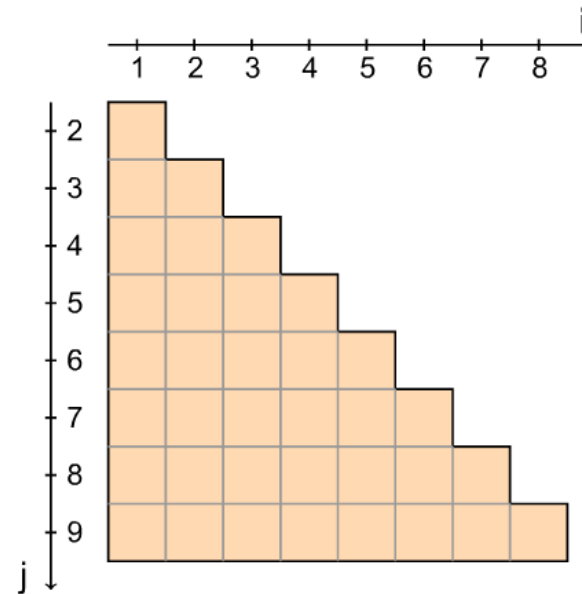
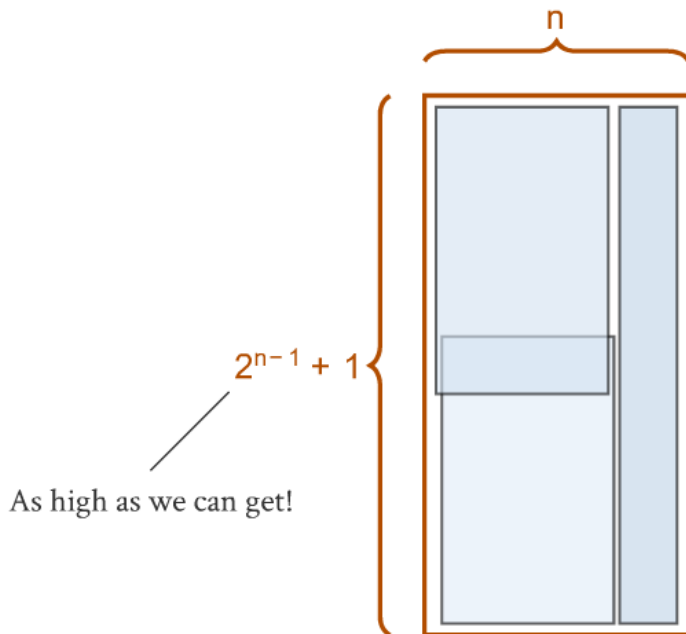
$$m^* \geq \Omega(\sqrt{2}^n) = \Omega(1.4142^n)$$

# A relaxation

$\Pi \in \{0, 1\}^{m \times n}$  is **quasi-Order-Regular**

iff  $\forall 0 < i < j < m, \exists 0 < k \leq n$  s.t.:

$$\Pi_{i,k} \neq \Pi_{i+1,k} = \Pi_{j,k} \neq \Pi_{j+1,k}$$



We can build quasi-Order-Regular matrices with  $2^{n-1} + 1$  rows!

# The state of the art and some inspiring ideas

**Upper bounds on  $m^*$**

$$m^* \leq O(2^n/n) \leq O(2^n)$$

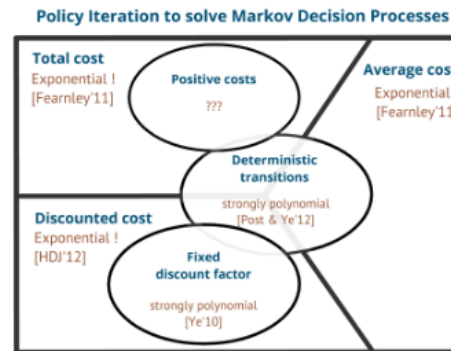
**Conjecture on  $m^*$**

$$m^* \sim O(\phi^n) = O(1.618^n)$$

**Lower bounds on  $m^*$**

$$m^* \geq \Omega(\sqrt{2}^n) = \Omega(1.4142^n)$$

# PART 1 : Exponential lower bounds for Policy Iteration



# PART 2 : Can we obtain good upper bounds?

The state of the art and some inspiring ideas

Upper bounds on  $m^*$        $m^* \leq O(2^n/n) \leq O(2^n)$

Conjecture on  $m^*$        $m^* \sim \Theta(\phi^n) = O(1.618^n)$

Lower bounds on  $m^*$        $m^* \geq \Omega(\sqrt{2}^n) = \Omega(1.4142^n)$

# Summary on the complexity of Policy Iteration

## Lower bounds

In general: **exponential** in the worst case

With some particular structures: **strongly polynomial**

What about the **average complexity**?

## Upper bounds

Best known bound:  $O(2^n/n)$

Our hope:  $O(1.618^n)$