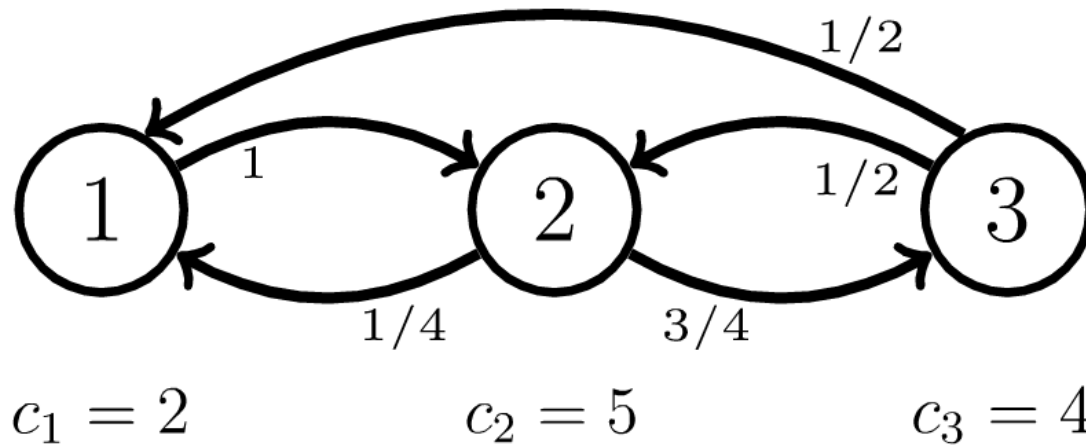# CDC'12

**The complexity of Policy Iteration is exponential for discounted Markov Decision Processes**

**Romain Hollanders**

Joint work with Raphaël Jungers and
Jean-Charles-Delvenne

# CDC'12

**The complexity of Policy Iteration is exponential for discounted Markov Decision Processes**

**Romain Hollanders**

Joint work with Raphaël Jungers and
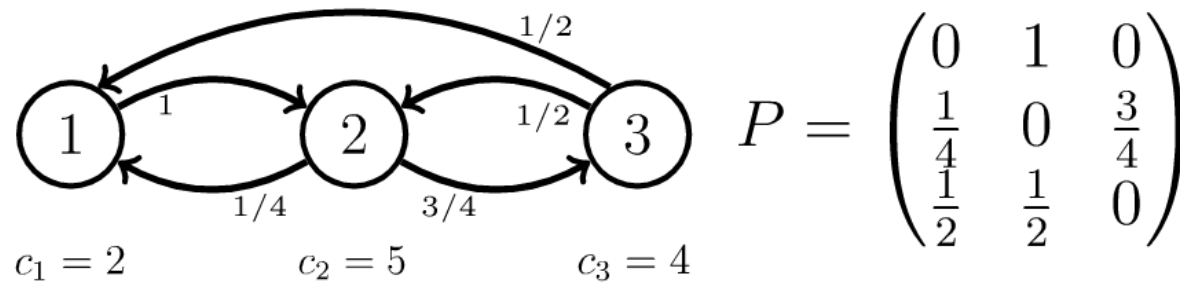Jean-Charles-Delvenne

# Markov Chains

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

$c_1 = 2$          $c_2 = 5$          $c_3 = 4$

$n_1^T = n_2^T \, P^k$

How much will I pay

# Markov Chains



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

$c_1 = 2 \qquad c_2 = 5 \qquad c_3 = 4$

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

$p_0^T = \begin{matrix} 1 & 0 & 0 \end{matrix}$

$p_1^T = \begin{matrix} 0 & 1 & 0 \end{matrix}$

$p_2^T = \begin{matrix} \frac{3}{4} & 0 & \frac{1}{4} \end{matrix}$

$p_3^T = \begin{matrix} \frac{1}{8} & \frac{7}{8} & 0 \end{matrix}$

$p_4^T = \begin{matrix} \frac{21}{32} & \frac{4}{32} & \frac{7}{32} \end{matrix}$

$\vdots$

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

- Total cost

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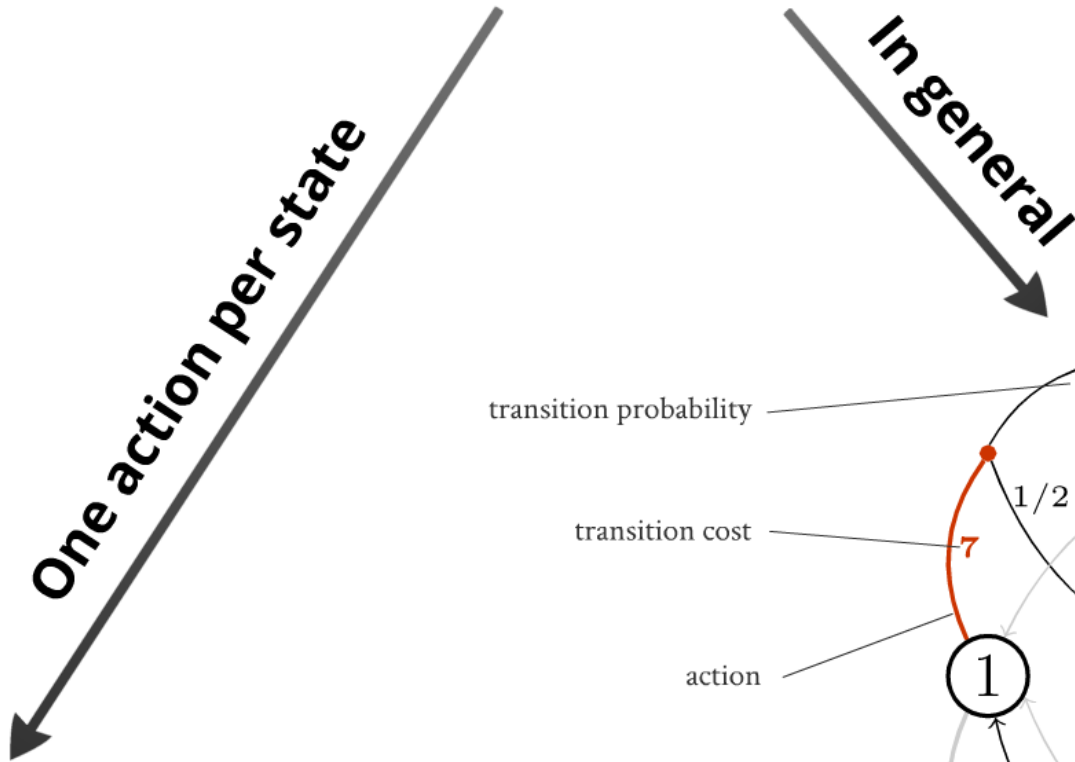
- Average cost

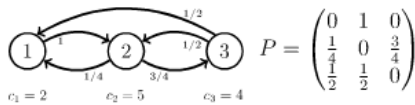$$x(1) = \lim_{H \to \infty} \frac{1}{H} \cdot \sum_{k=0}^{H} p_k^T c$$

- Discounted cost

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

# Markov Decision Processes

**One action per state**

**In general**

## Markov Chains



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

$c_1 = 2 \qquad c_2 = 5 \qquad c_3 = 4$

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

$p_0^T = 1 \quad 0 \quad 0$
$p_1^T = 0 \quad 1 \quad 0$
$p_2^T = \frac{3}{4} \quad 0 \quad \frac{1}{4}$
$p_3^T = \frac{1}{8} \quad \frac{7}{8} \quad 0$
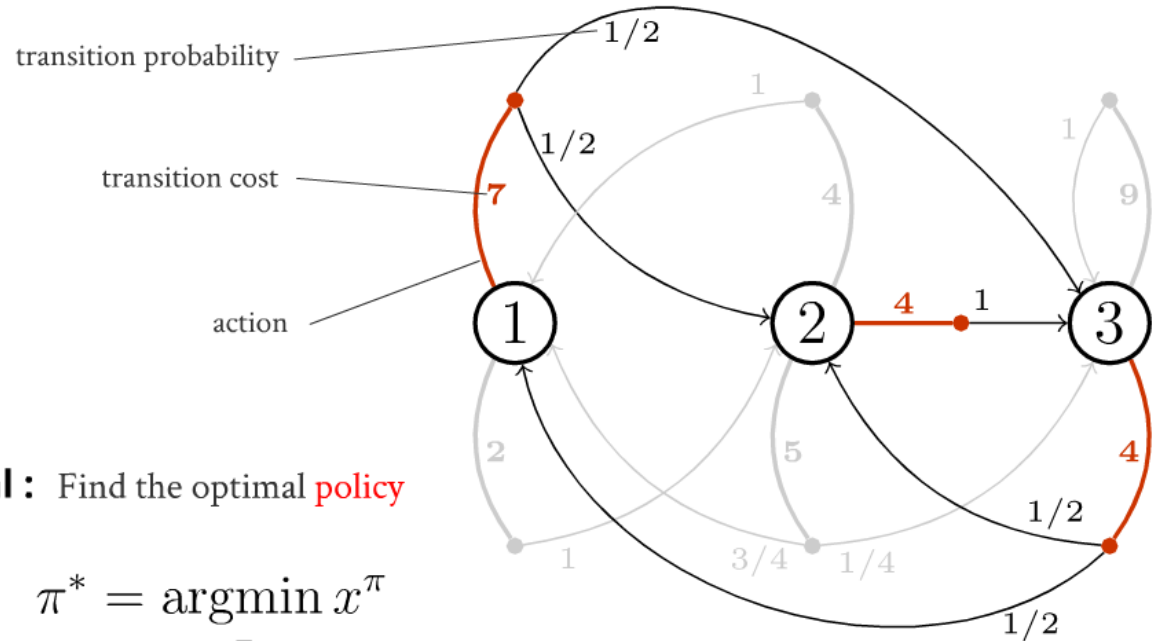$p_4^T = \frac{21}{32} \quad \frac{4}{32} \quad \frac{7}{32}$
$\vdots$

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

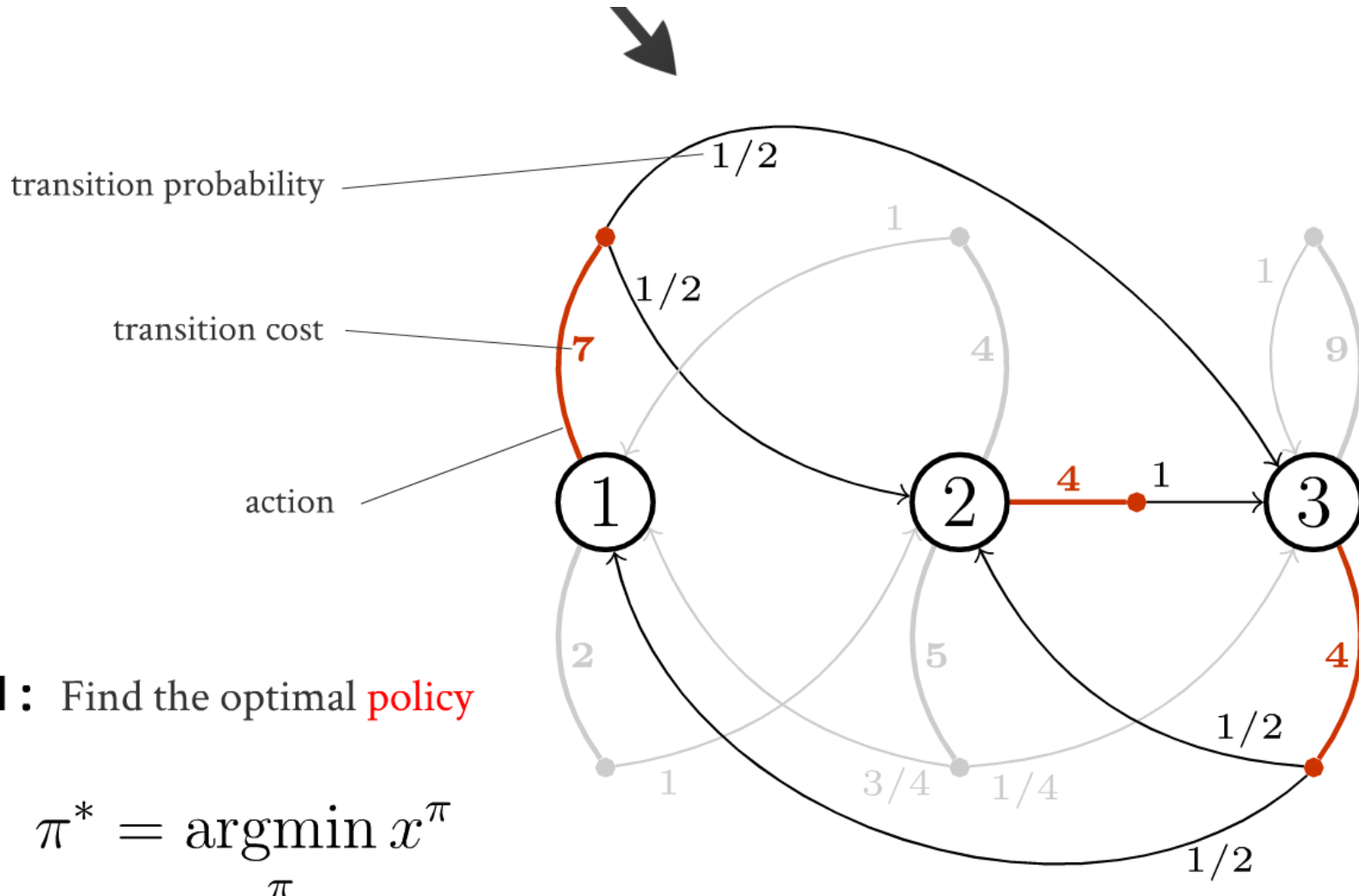● Discounted cost
$$x(1) = \sum_{k=0}^{\infty} \gamma^k p_k^T c$$

transition probability

transition cost

action



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

$$\pi^* = \underset{\pi}{\arg\min} \, x^\pi$$

The answer depends on the chosen objective function :

● Total cost

● Average cost

● Discounted cost

transition probability

transition cost

action

**Goal :** Find the optimal policy

$$\pi^* = \underset{\pi}{\mathrm{argmin}}\; x^\pi$$

The answer depends on the chosen objective function :

- Total cost
- Average cost
- Discounted cost

**Markov Decision Processes**

One action per state

In general

Resolution

**Linear Programming**

**Value Iteration**

**Policy Iteration**

transition probability

transition cost

action

**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 Chains**

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 $\pi_0$

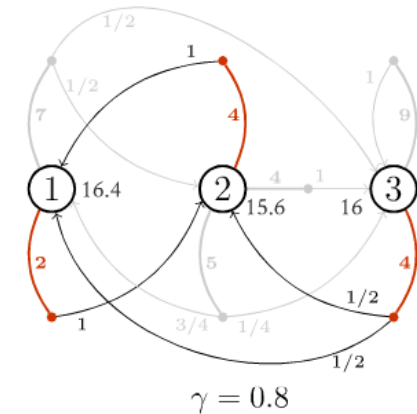while  $\pi_k \neq \pi_{k-1}$

    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



$\gamma = 0.8$

# Policy Iteration

0. Choose an initial policy $\pi_0$

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

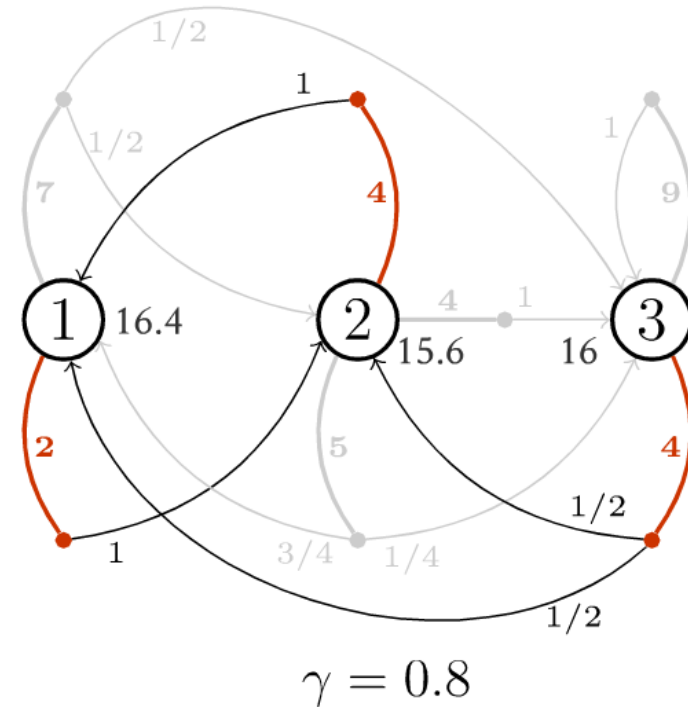    1. Evaluate $\pi_k$ (Bellman)

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

    2. Improve $\pi_k$

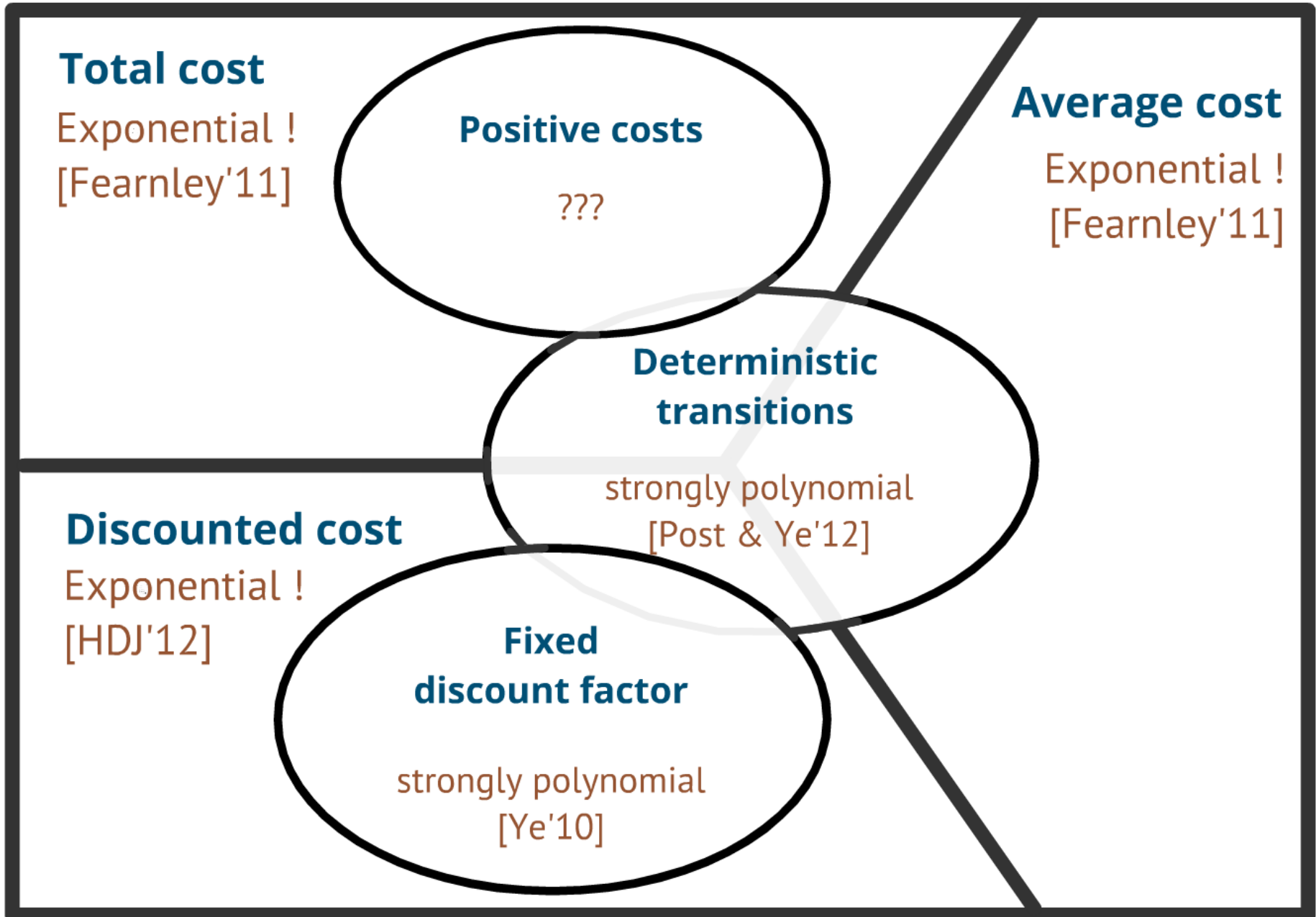$$\pi_{k+1} = \operatorname*{argmin}_{\pi} \, c^{\pi} + \gamma \, P^{\pi} x^{\pi_k}$$
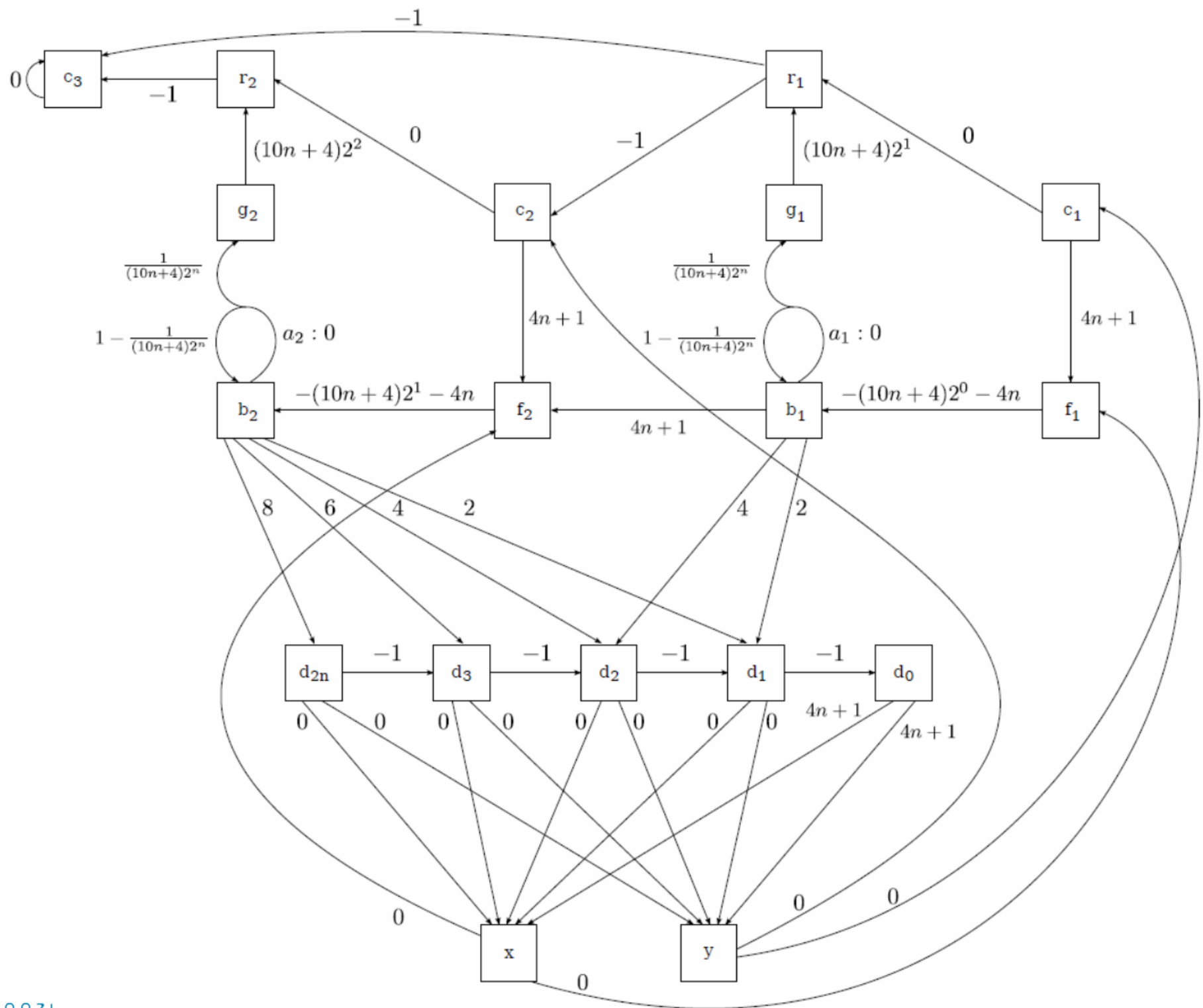
$$k \; \leftarrow \; k+1$$

end



$\gamma = 0.8$

# Policy Iteration to solve Markov Decision Processes

**Total cost**

Exponential !
[Fearnley'11]

**Positive costs**

???

**Average cost**

Exponential !
[Fearnley'11]

**Deterministic transitions**

strongly polynomial
[Post & Ye'12]

**Discounted cost**

Exponential !
[HDJ'12]

**Fixed discount factor**

strongly polynomial
[Ye'10]

# Policy Iteration to solve Markov Decision Processes

**Total cost**

Exponential !
[Fearnley'11]

**Positive costs**

???

**Average cost**

Exponential !
[Fearnley'11]

**Deterministic
transitions**

strongly polynomial
[Post & Ye'12]

**Discounted cost**

Exponential !
[HDJ'12]

**Fixed
discount factor**

strongly polynomial
[Ye'10]

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

➡ How much perturbation?

$x^{\pi'}(s)$

$\tilde{x}^{\pi'}(s)$ $\left.\begin{array}{c}\end{array}\right\} \leq F(n, \delta, \kappa)\,\varepsilon$

$\tilde{x}^{\pi_k}(s)$

$x^{\pi_k}(s)$ $\leq F(n, \delta, \kappa)\,\varepsilon$

$\geq G(n, \delta)$

➡ OK for some $\varepsilon \sim \dfrac{1}{2^{q(n,\delta,\kappa)}}$

# Policy Iteration to solve Markov Decision Processes

**Total cost**

Exponential !
[Fearnley'11]

**Positive costs**

???

**Average cost**

Exponential !
[Fearnley'11]

**Deterministic transitions**

strongly polynomial
[Post & Ye'12]

**Discounted cost**

Exponential !
[HDJ'12]

**Fixed discount factor**

strongly polynomial
[Ye'10]