

Solving Strong Stackelberg Equilibrium in Stochastic Games

Víctor Bucarey López

vbucarey@ing.uchile.cl

INRIA – LILLE

December 2017

Joint Collaboration...

- Fernando Ordoñez
Departamento Ingeniería Industrial Universidad de Chile.
- Eugenio Della Vecchia
FCEIA - Universidad Nacional de Rosario.
- Alain Jean Marie
INRIA - Research center of Sophia-Antipolis Méditerranée

Project DyGaMe STIC AmSud

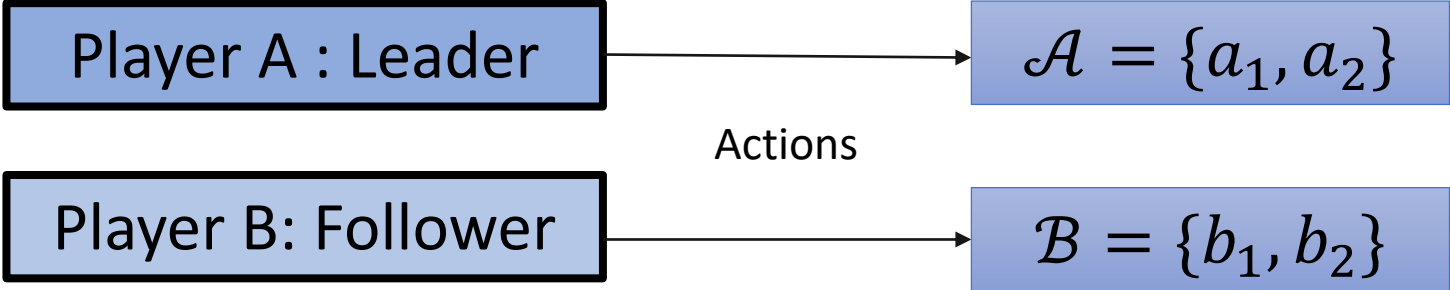
Stackelberg Game



STRONG STACKELBERG EQUILIBRIUM

- Leader commits to a payoff maximizing strategy.
 - Follower best responds
 - Follower breaks ties in favor of the leader

Example



		b_1	b_2
x	a_1	$(10, -10)$	$(-5, 6)$
$1-x$	a_2	$(-8, 4)$	$(6, -4)$

If B plays b_1 leader's reward will be $10x + -8(1-x)$

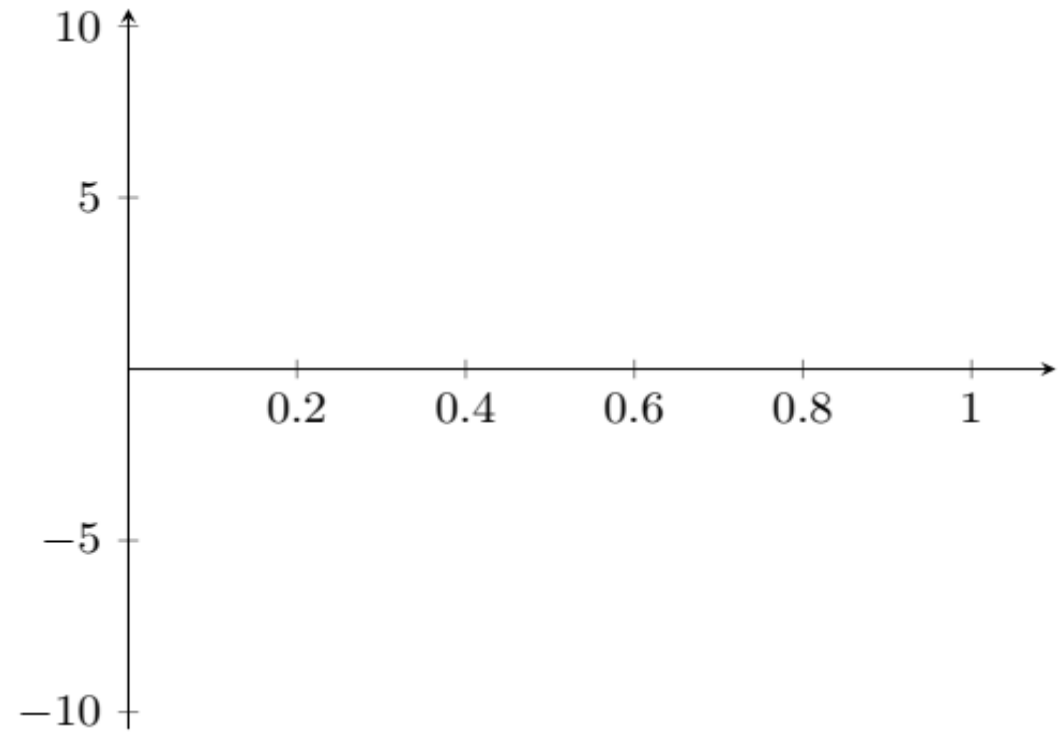
If B plays b_1 his expected reward will be $-10x + 4(1-x)$

If B plays b_2 leader's reward will be $-5x + 6(1-x)$

If B plays b_2 his expected reward will be $6x + -4(1-x)$

Example

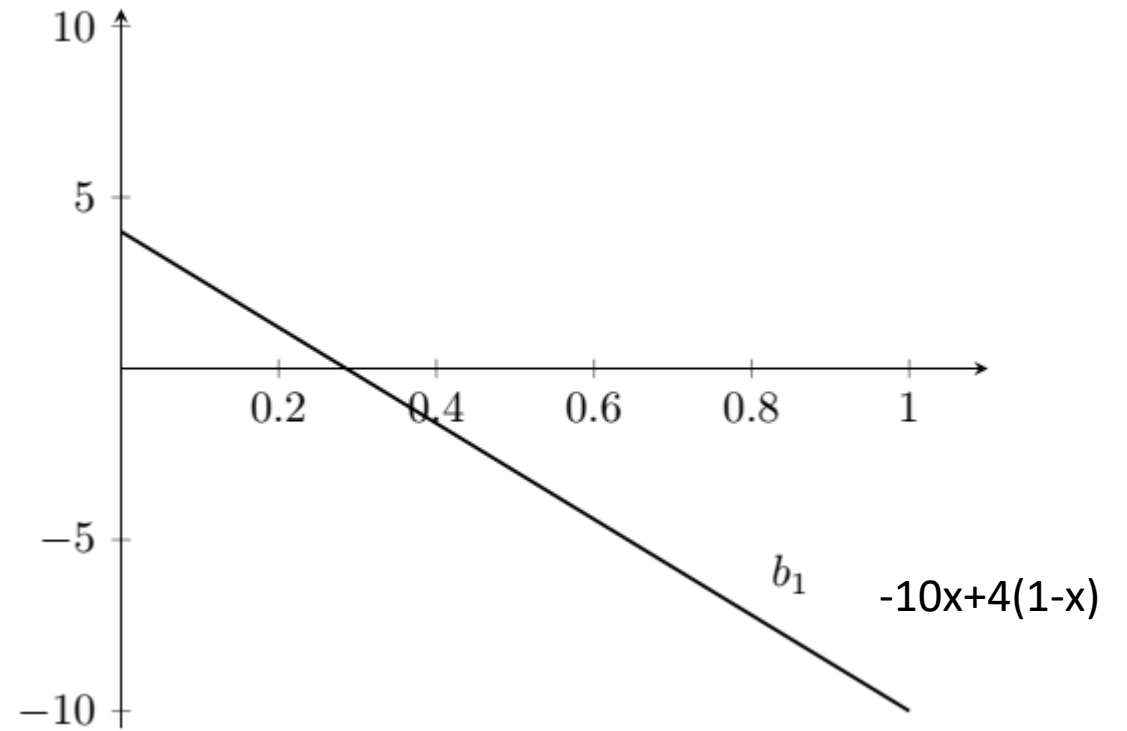
	b_1	b_2
x	$(10, -10)$	$(-5, 6)$
$1-x$	$(-8, 4)$	$(6, -4)$



Follower

Example

	b_1	b_2
x	$(10, -10)$	$(-5, 6)$
$1-x$	$(-8, 4)$	$(6, -4)$



Follower

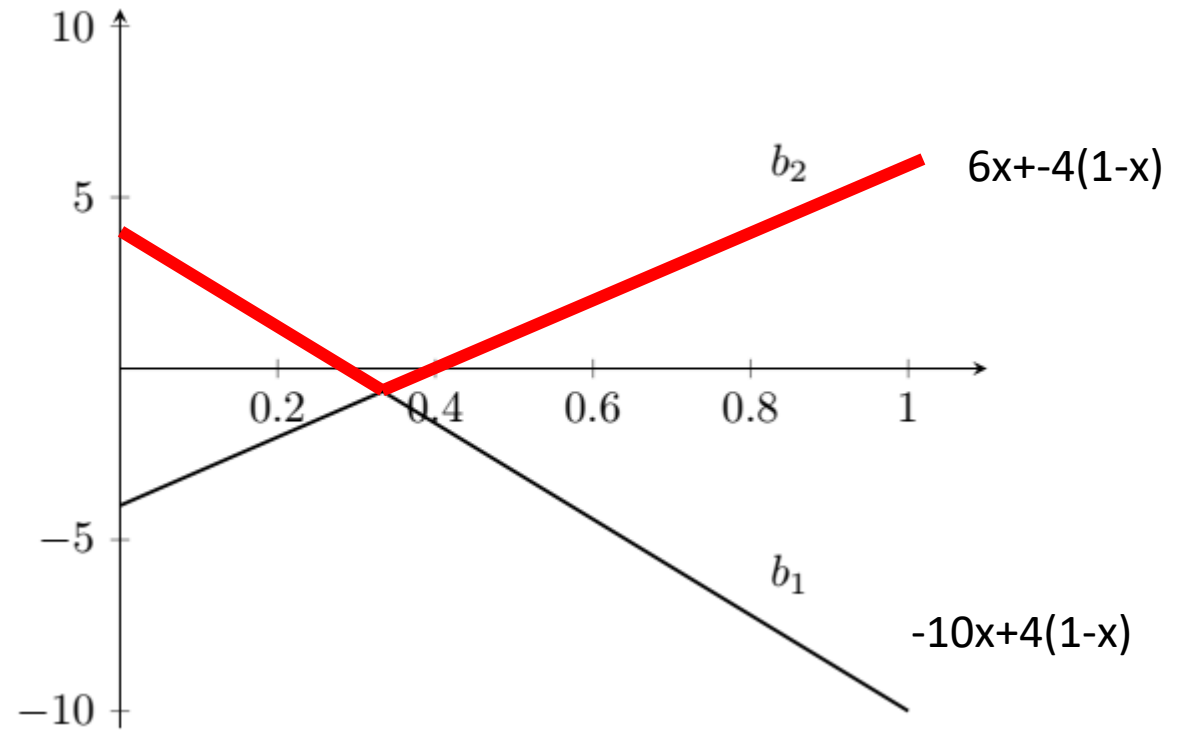
Example

	b_1	b_2
x	$(10, -10)$	$(-5, 6)$
$1-x$	$(-8, 4)$	$(6, -4)$

Best Response:

$$g(x) = \begin{cases} b_2 & \text{if } x > \frac{1}{3} \\ b_1 & \text{if } x < \frac{1}{3} \end{cases}$$

$$x = 1/3 ?$$



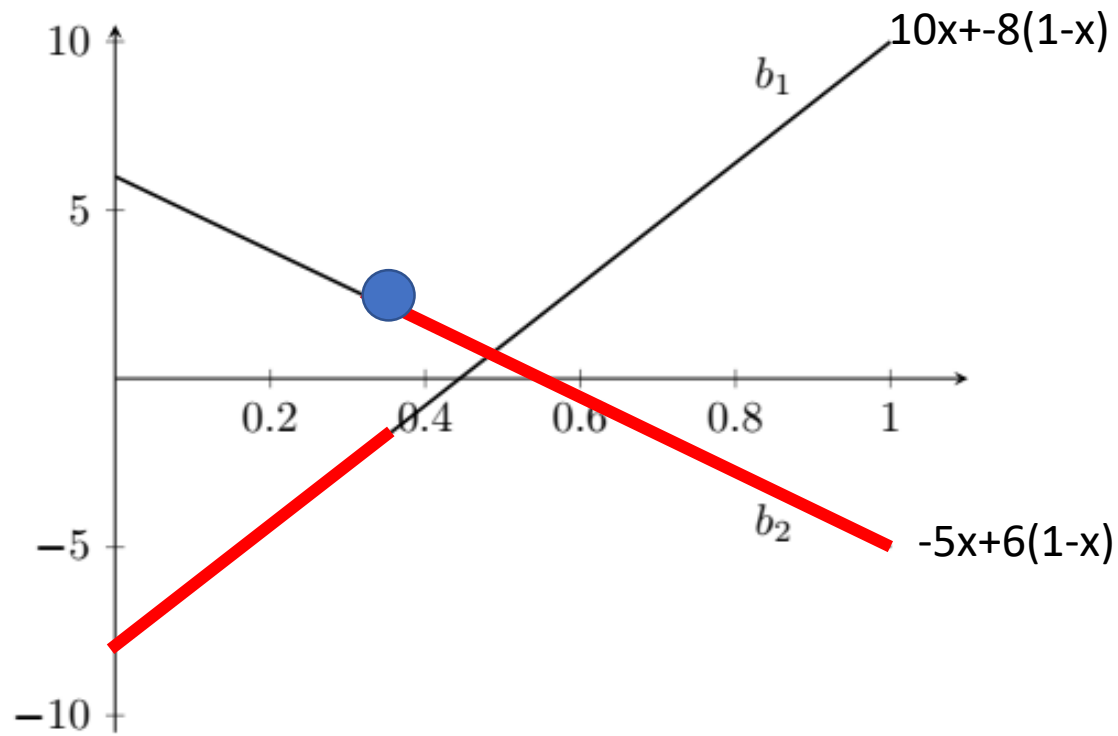
Follower

Example

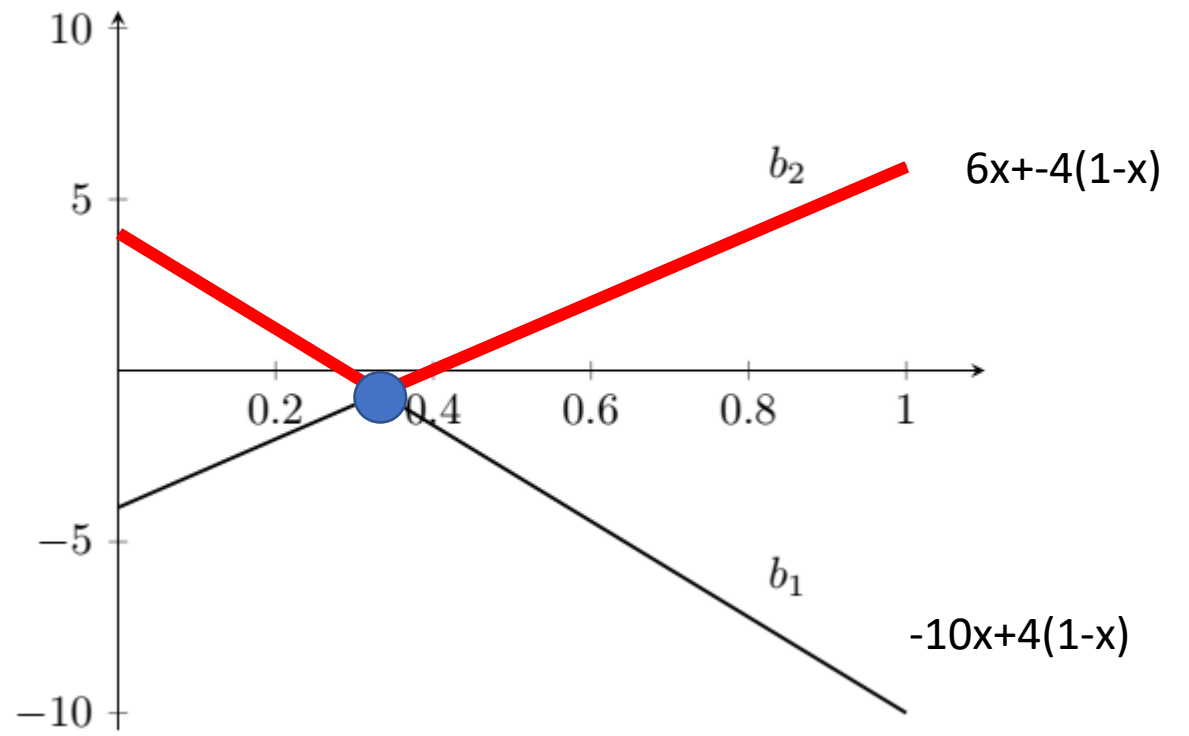
x

1-x

	b_1	b_2
a_1	(10, -10)	(-5, 6)
a_2	(-8, 4)	(6, -4)



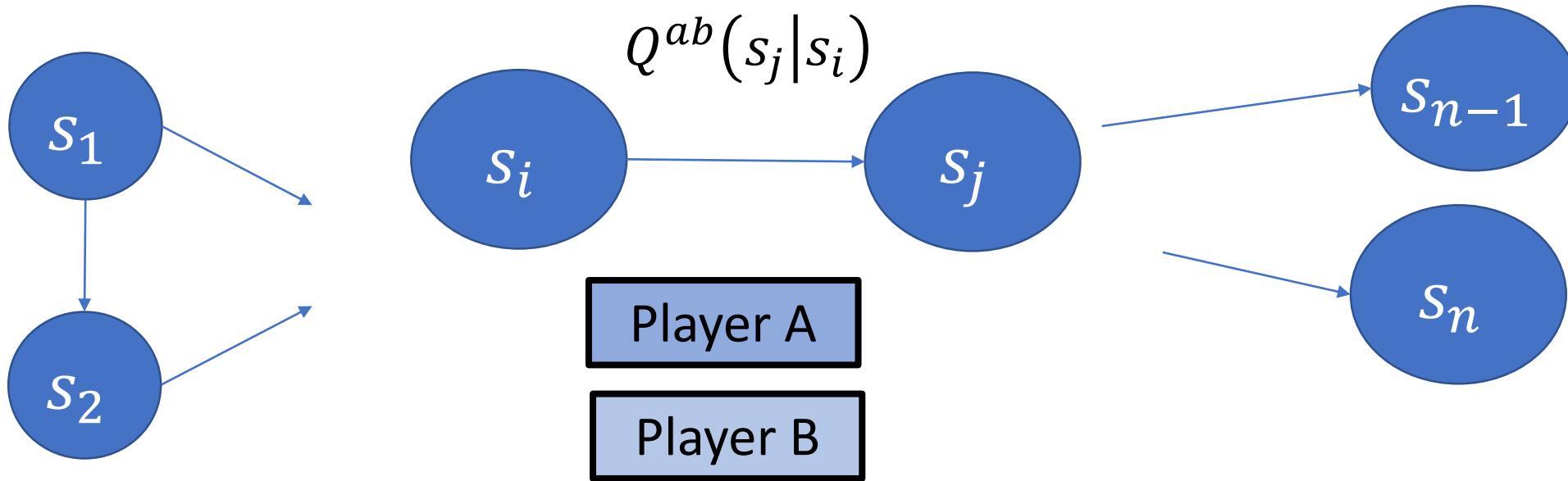
Leader



Follower

Stochastic Games

$$\mathcal{G} = (\mathcal{S}, \mathcal{A}, \mathcal{B}, r_A, r_B, Q, \beta_A, \beta_B, \tau)$$



Payoffs:

$$r_A = \{r_A^{ab}(s_i)\}$$
$$r_B = \{r_B^{ab}(s_i)\}$$

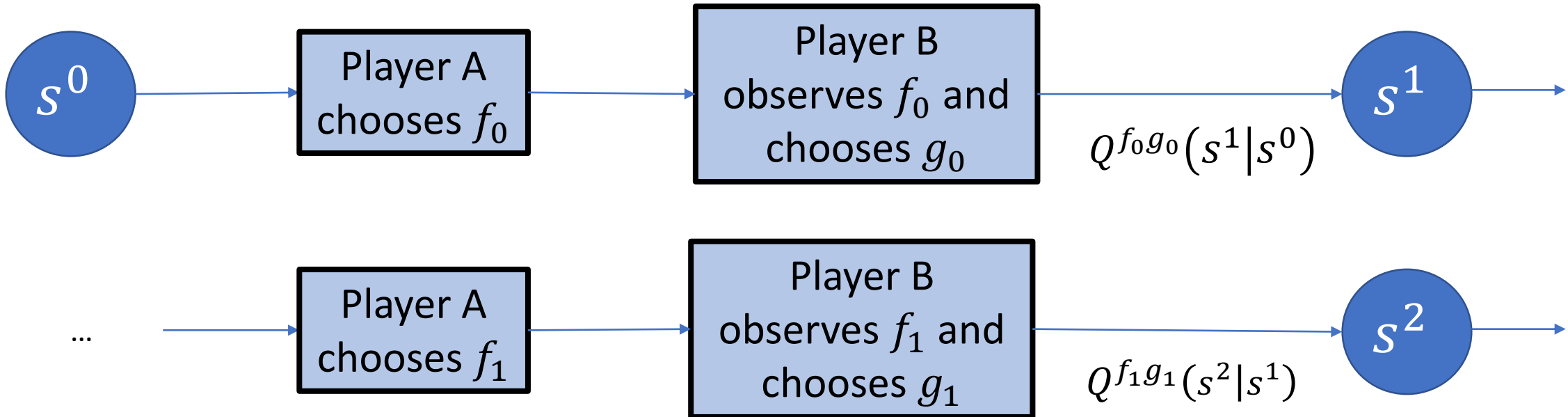
Discount Factors:

$$\beta_A, \beta_B \in [0, 1)$$

Time horizon

$$\tau \in \mathbb{N} \cup \{+\infty\}$$

Stochastic Games



Feedback Policies:
 $\pi = \pi(s, t)$
 $= \{f_1, \dots, f_\tau\}$

Stationary Policies:
 $\pi = \pi(s)$
 $= \{f, \dots, f\}$

Framework

- **General Objectives**

- Existence and characterization of value functions
- Existence of equilibrium strategies
- Algorithms to compute them

- **State of the art**

- For finite horizon, Stackelberg equilibrium in stochastic games via Dynamic programming.
- Mathematical programming approach to compute stationary values.

Our contribution

- We define suitable Dynamic Programming operators.
- We used it to characterize value functions and to prove existence and unicity of stationary values forming a Strong Stackelberg Equilibrium for a family of problems.
- We define Value Iteration and Policy Iteration for this family and prove its convergence.
- We prove via counterexample that this methodology is not always applicable for the general case.

Stackelberg Equilibrium in Stochastic Games

(π, γ)



Value functions

$$v_A^{\pi, \gamma}(s) = \mathbb{E}^{\pi, \gamma}(s) \left[\sum_{t=0}^{\tau} \beta_A^t r_A^{A_t B_t}(S_t) \right]$$

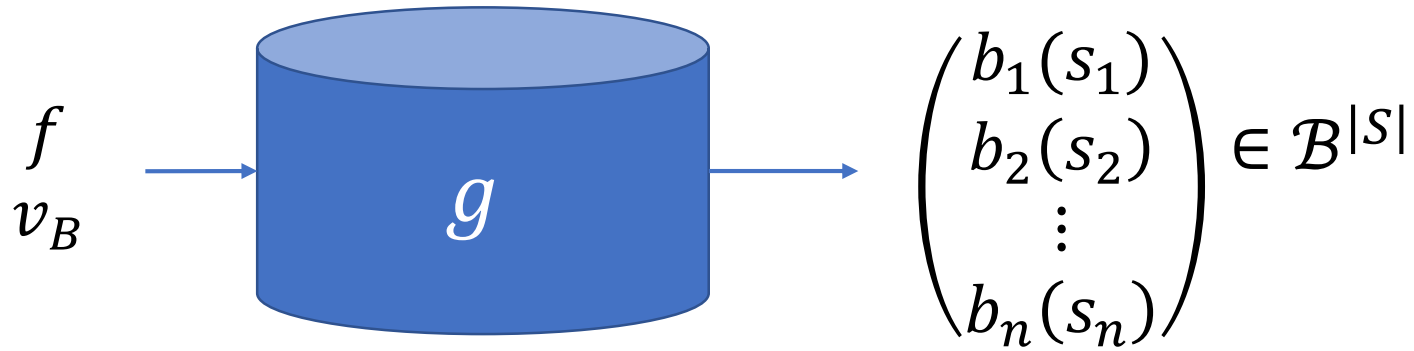
$$v_B^{\pi, \gamma}(s) = \mathbb{E}^{\pi, \gamma}(s) \left[\sum_{t=0}^{\tau} \beta_B^t r_B^{A_t B_t}(S_t) \right]$$

Stackelberg Equilibrium:

(π^*, γ^*)

$$v_A^{\pi^*, \gamma^*}(s) = \max_{\pi} v_A^{\pi, \gamma^*}(s)$$
$$\gamma^* \in \operatorname{argmax}_{\gamma} v_B^{\pi^*, \gamma}(s)$$

Best response functional



Given a stationary policy and future values, g computes the best actions to perform in each state.

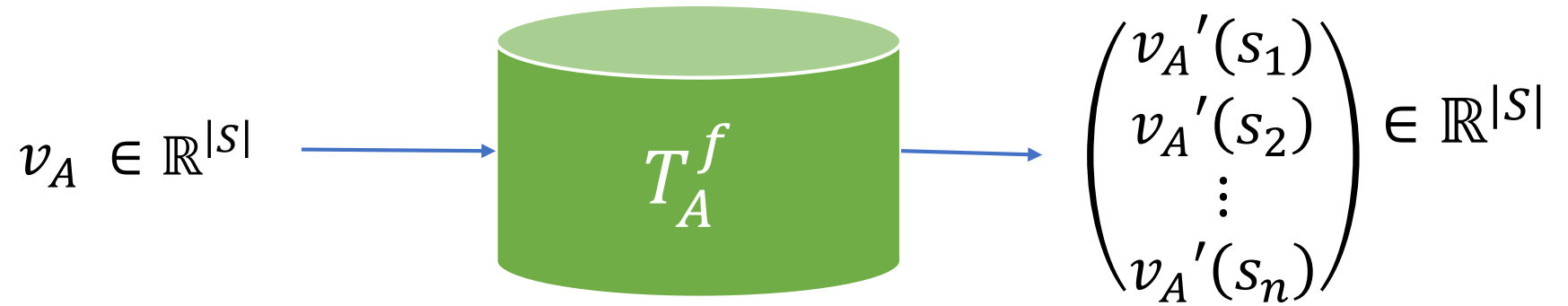
$$g(f, v_B)(s) = \operatorname{argmax}_{b \in \mathcal{B}} \sum_{a \in \mathcal{A}} f(a, s) \left[r_B^{ab}(s) + \beta_B \sum_{z \in S} Q^{ab}(z|s) v_B(z) \right]$$



- **Myopic follower strategies (MFS)**

$$g(f, v_B) = g(f)$$

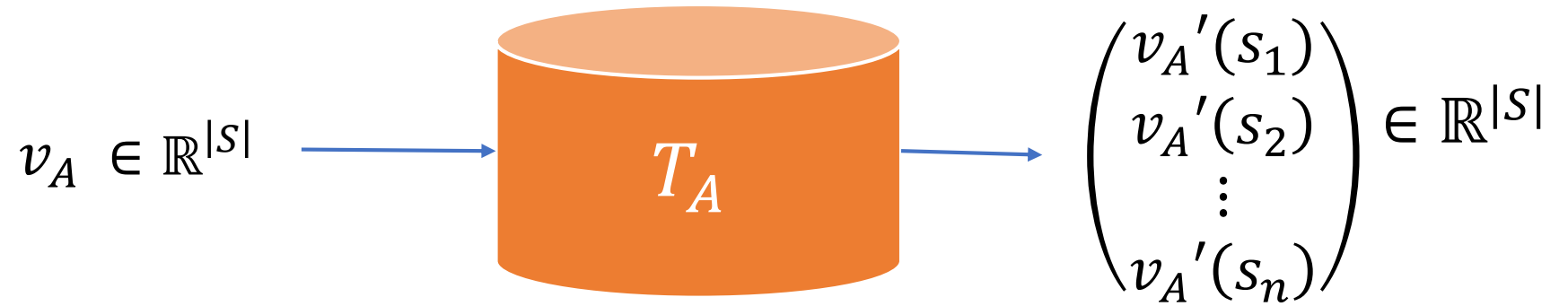
MFS case



For a fixed stationary policy f , T_A^f computes the value of applying this policy.

$$T_A^f(v_A)(s) = \sum_{a \in \mathcal{A}} f(a, s) \left[r_B^{ag(f)(s)}(s) + \beta_B \sum_{z \in S} Q^{ag(f)(s)}(z|s) v_A(z) \right]$$

MFS case



T_A computes the value of the best policy for the leader.

$$T_A(v_A)(s) = \max_{f \in \mathcal{P}(\mathcal{A})} T_A^f(v_A)(s)$$

MFS case



- **Theorem 1.**

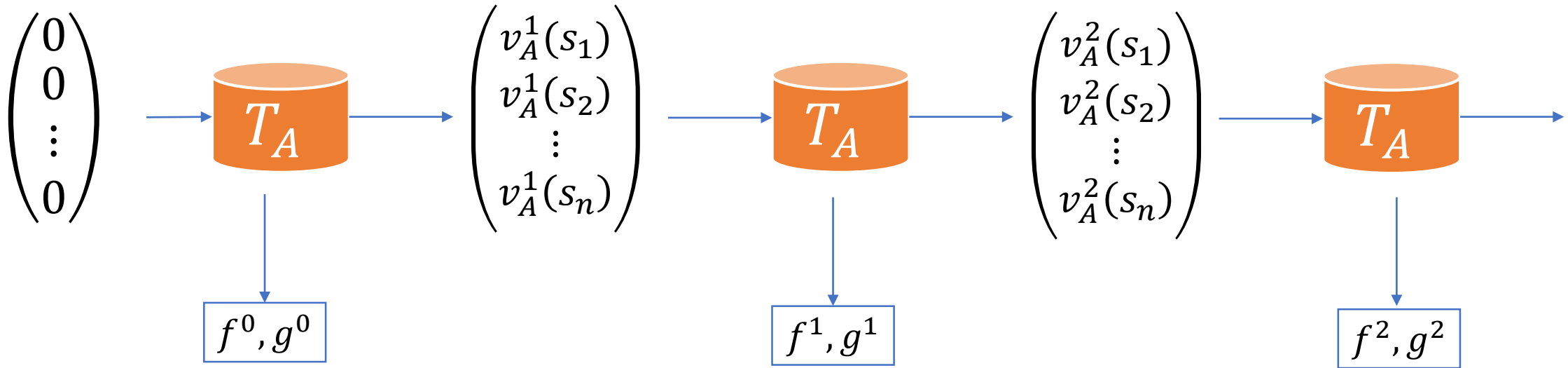
- T_A^f, T_A are monotone.
- For any stationary strategy f , the operator T_A^f is a contraction on $(\mathbb{R}^{|S|}, \|\cdot\|_\infty)$ of modulus β_A .
- The operator T_A is a contraction on $(\mathbb{R}^{|S|}, \|\cdot\|_\infty)$ of modulus β_A .

- **Theorem 2.**

There exists a equilibrium value function v_A^* and it is the unique solution of $v_A^* = T_A(v_A^*)$. Moreover, the pair f^* and $g(f^*)$ which maximizes the RHS of (1) are the equilibrium strategies.



Value Iteration algorithm



Repeat until
convergence

Theorem 3.

The sequence of value functions v_A^n converges to v_A^* . Furthermore, v_A^* is the fixed point of T_A with the following bound:

$$\|v_A^* - v_A^n\| \leq \beta_A^n \frac{\|r_A\|_\infty}{1 - \beta_A}$$



Value Iteration algorithm

Algorithm 1 Value function iteration: Infinite horizon

Require: $\varepsilon > 0$

1: Initialize with $n = 1$, $v_A^0(s) = 0$ for every $s \in \mathcal{S}$ and $v_A^1 = T_A(v_A^0)$

2: **while** $\|v_A^n - v_A^{n-1}\|_\infty > \varepsilon$ **do**

3: Compute v_A^{n+1} by

$$v_A^{n+1}(s) = T_A(v_A^n)(s) .$$

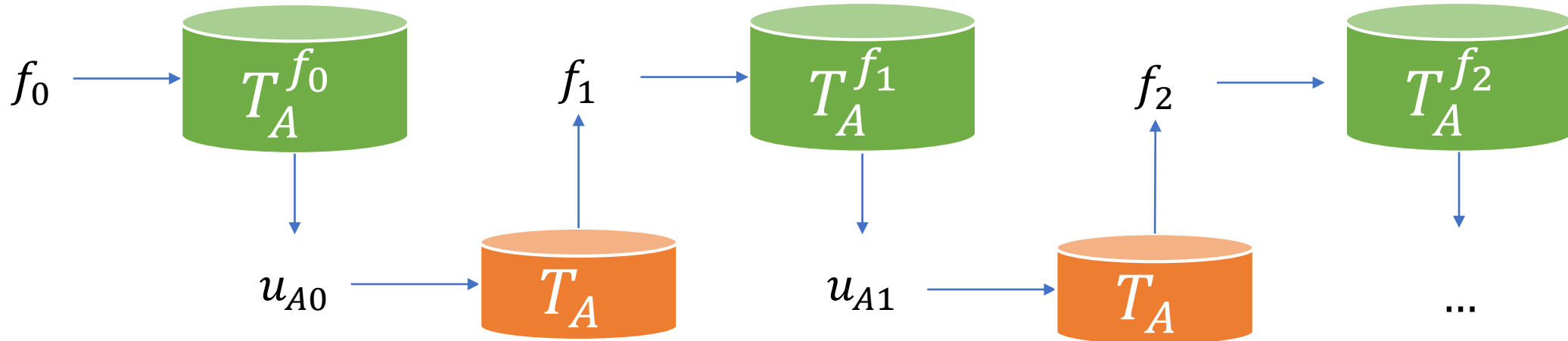
 Finding f^* and $g^*(f)$ at stage n .

4: $n := n + 1$

5: **end while**

6: **return** Stationary Stackelberg policies $\pi^* = \{f^*, \dots\}$ and $\gamma^* = \{g^*, \dots\}$

Policy Iteration algorithm



Repeat until
convergence

Theorem 4.

The sequence of functions $u_{A,n}$ verifies $u_{A,n} \uparrow v_A^*$. Furthermore, if for any $n \in \mathbb{N}$, $u_{A,n} = u_{A,n+1}$ then it is true that $u_{A,n} = v_A^*$.



Policy Iteration algorithm

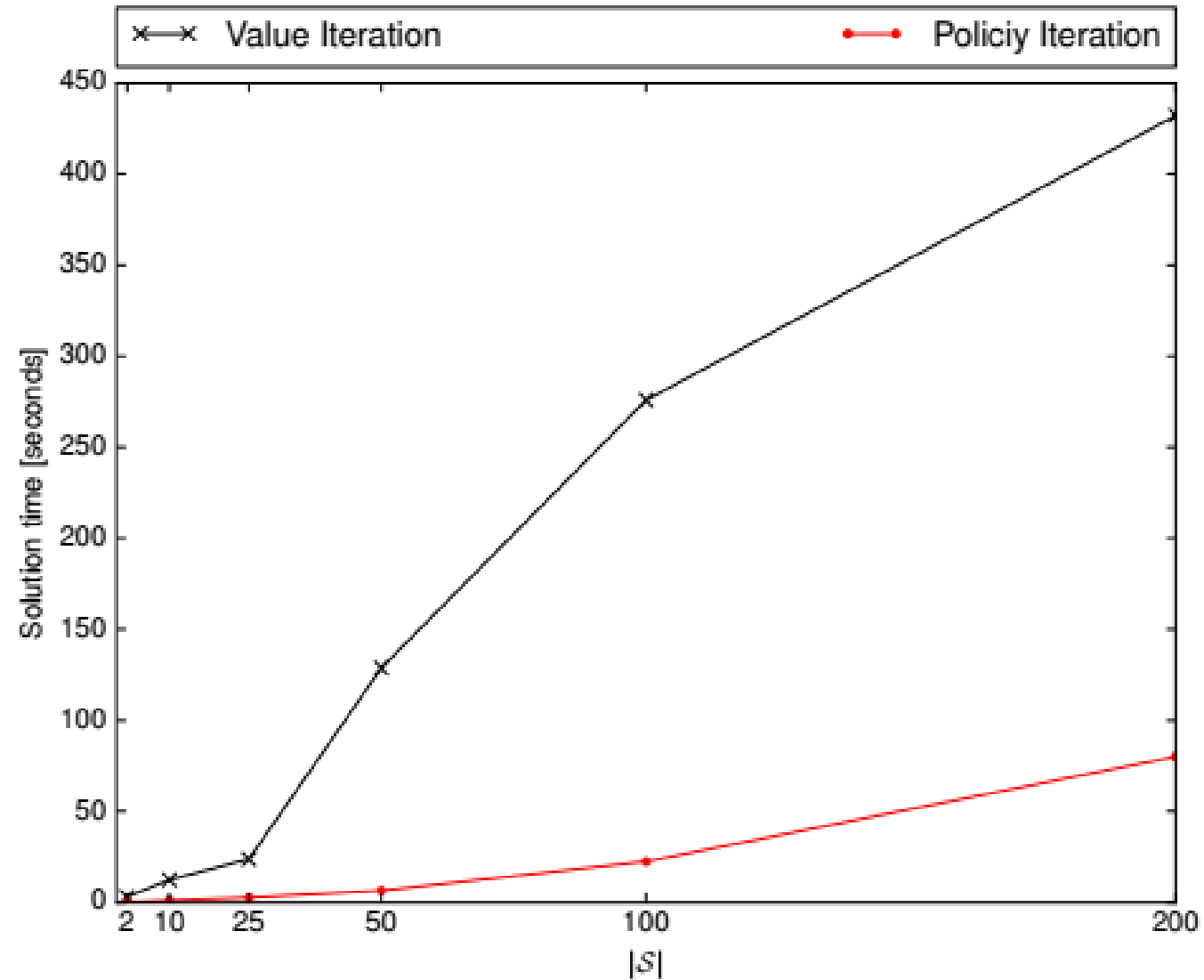
Algorithm 2 Policy Iteration (PI)

- 1: Choose a stationary Stackelberg pair $(f_0, g(f_0))$.
- 2: **while** $\|u_{A,n} - u_{A,n+1}\| > \varepsilon$ **do**
- 3: Evaluation Phase: Find $u_{A,n}$ fixed point of the operator $T_A^{f_n}$.
- 4: Improvement Phase: Find a strategy f_{n+1} such that

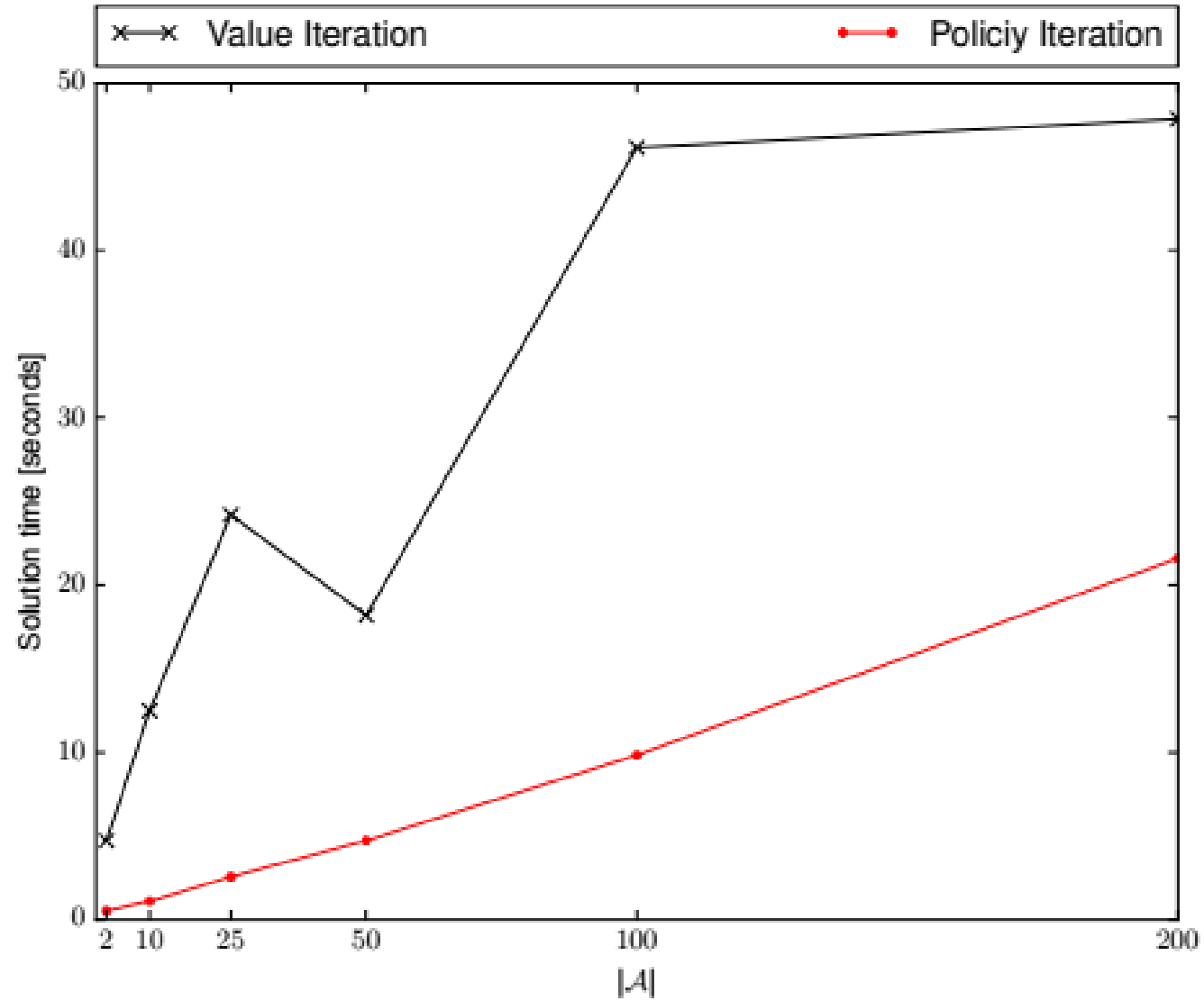
$$T_A^{f_{n+1}}(u_{A,n}) = T_A(u_{A,n}) .$$

- 5: $n := n + 1$
 - 6: **end while**
 - 7: **return** Stationary Stackelberg policies $\pi^* = \{f^*, \dots\}$ and $\gamma^* = \{g(f^*), \dots\}$
-

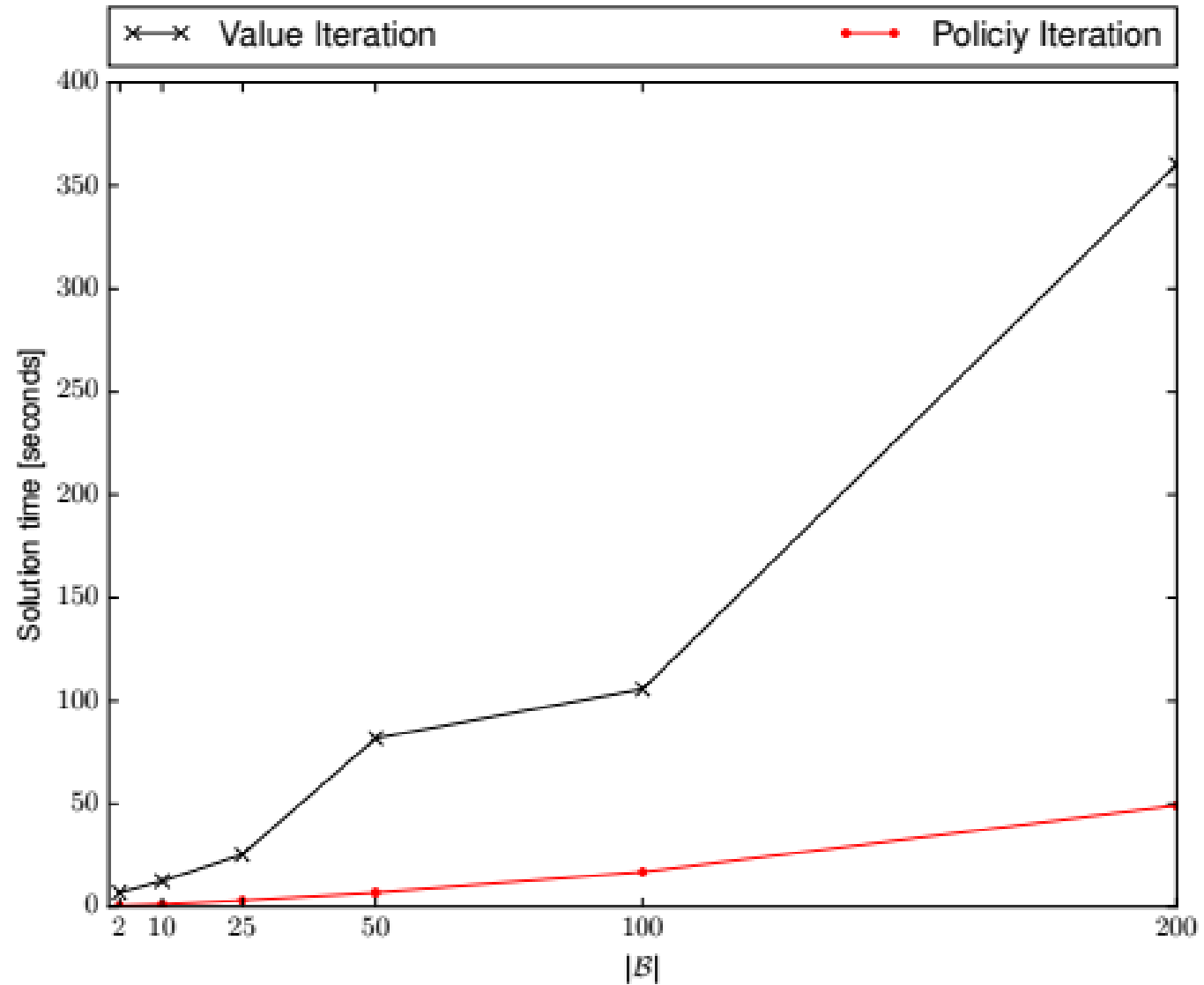
Computational Results



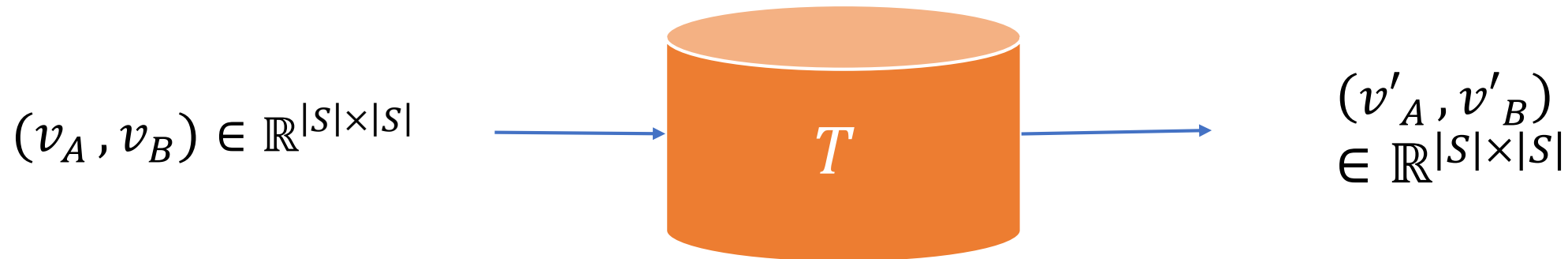
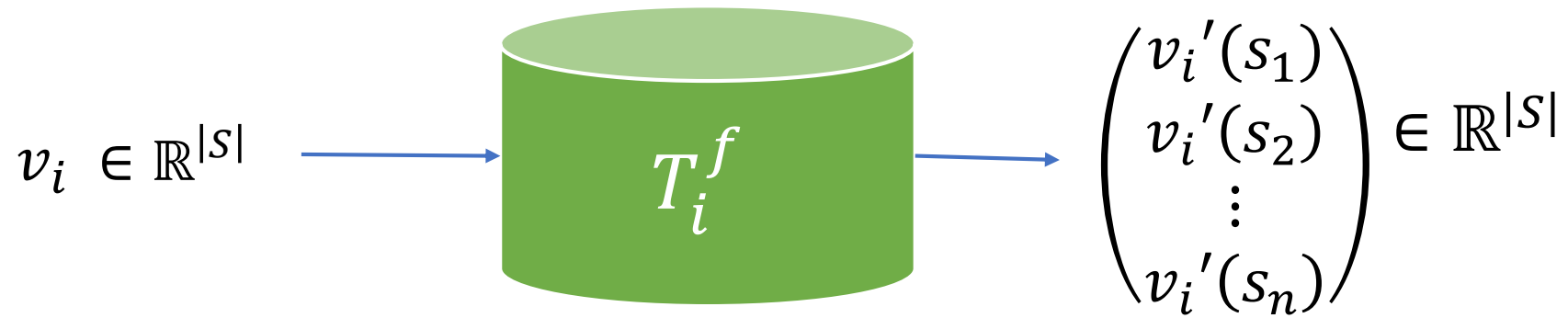
Computational Results



Computational Results



General case



$$(T(v_A, v_B))(s) = \left(\max_{f \in \mathcal{P}(\mathcal{A}_s)}, T_A^{f, g(f, v_B)}(v_A)(s), T_B^{f^*, g(f^*, v_B)}(v_B)(s) \right)$$

General case



Algorithm 3 Value Iteration (VI): Finite horizon for the general case

1: Initialize with $v_A^{\tau+1}(s) = v_B^{\tau+1}(s) = 0$ for every $s \in \mathcal{S}$

2: **for** $t = \tau, \dots, 0$, and for every $s \in \mathcal{S}$ **do**

3: Solve

$$(v_A^t(s), v_B^t(s)) = T(v_A^{t+1}, v_B^{t+1})(s) \quad \forall s \in \mathcal{S}$$

 Finding f_t^* and g_t^* SSE strategies at stage t .

4: **end for**

5: **return** Stackelberg policies $\pi^* = \{f_0^*, \dots, f_\tau^*\}$ and $\gamma^* = \{g_0^*, \dots, g_\tau^*\}$

This algorithm returns an Strong Stackelberg Equilibrium in **feedback policies** for the τ -horizon problem.

What about **stationary policies**?

Example



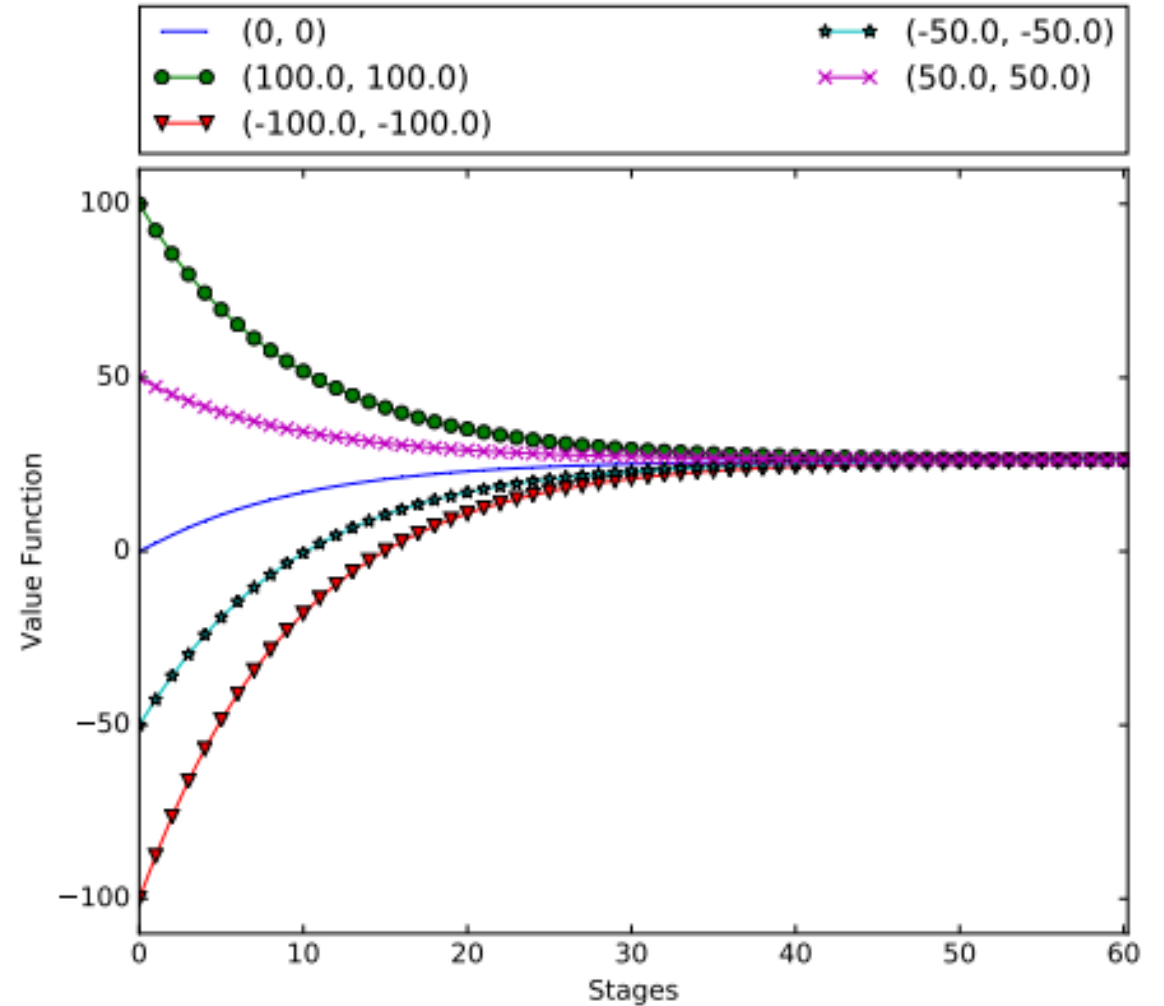
	b_1	b_2
a_1	$(\frac{1}{2}, \frac{1}{2})$ $(10, -10)$	$(0, 1)$ $(-5, 6)$
a_2	$(\frac{1}{4}, \frac{3}{4})$ $(-8, 4)$	$(1, 0)$ $(6, -4)$

State s_1

	b_1	b_2
a_1	$(\frac{1}{2}, \frac{1}{2})$ $(7, -5)$	$(0, 1)$ $(-1, 6)$
a_2	$(\frac{1}{4}, \frac{3}{4})$ $(-3, 10)$	$(1, 0)$ $(2, -10)$

State s_2

$$\beta_A = \beta_B = 0.9$$



Example

	b_1	b_2
a_1	(1, 0) / (1, -1)	(0, 1) / (0, 1)
a_2	(0, 1) / (-1, 1)	(0, 1) / (-1, -1)

State s_1

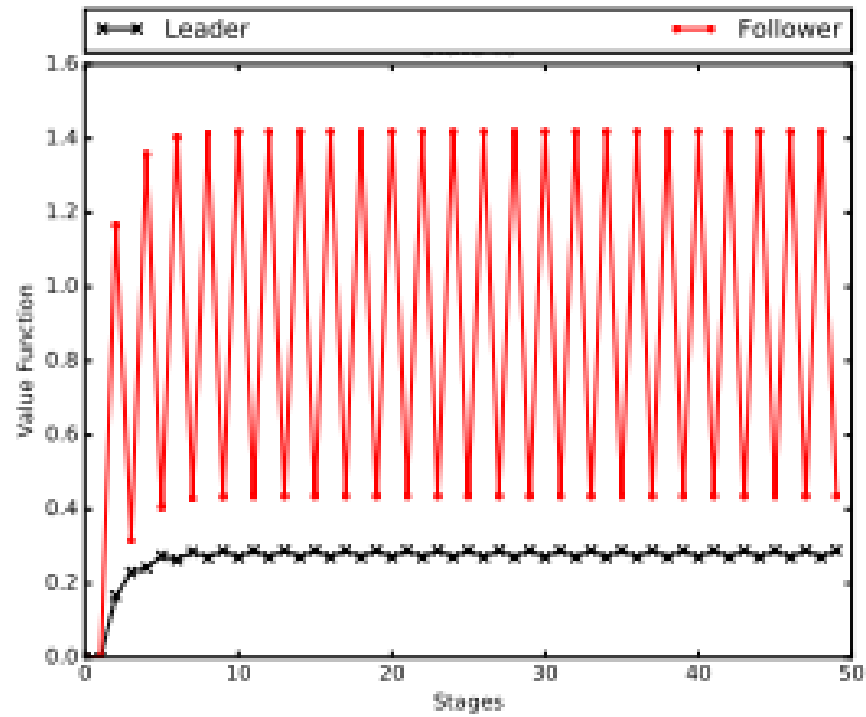
	b_1	b_2
a_1	(0, 1) / (-1, 0)	(1, 0) / (0, 1)
a_2	(1, 0) / (0, 1)	(0, 1) / (1, -1)

State s_2

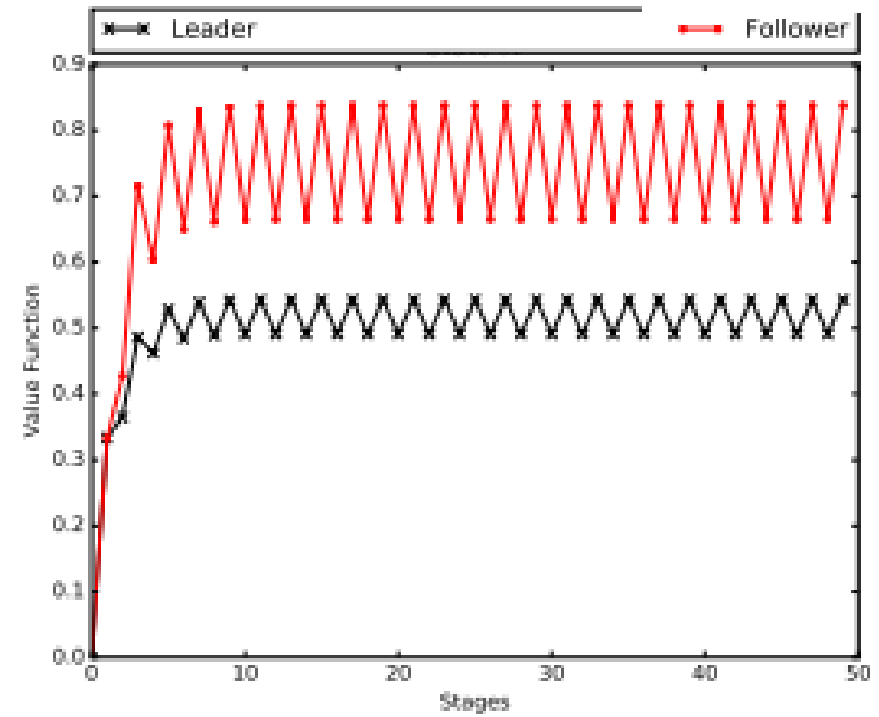


$$\beta_A, \beta_B = \frac{1}{2}$$

Example



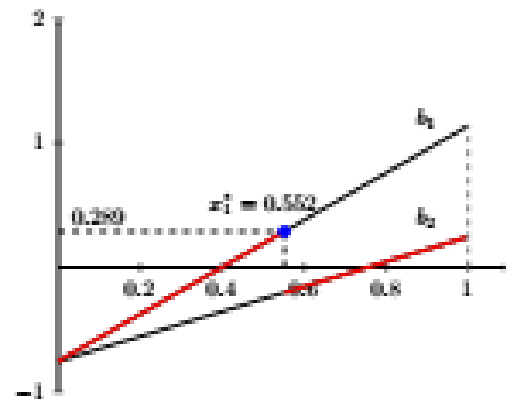
State s_1



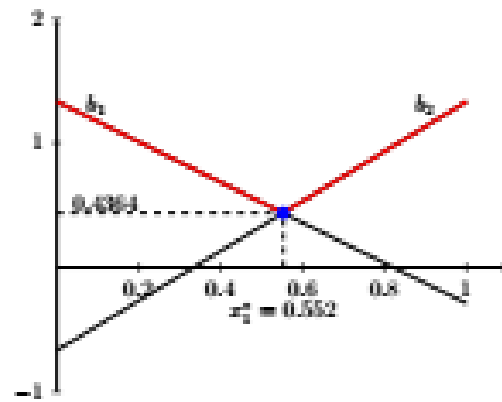
State s_2

Iteration 14

State s_1

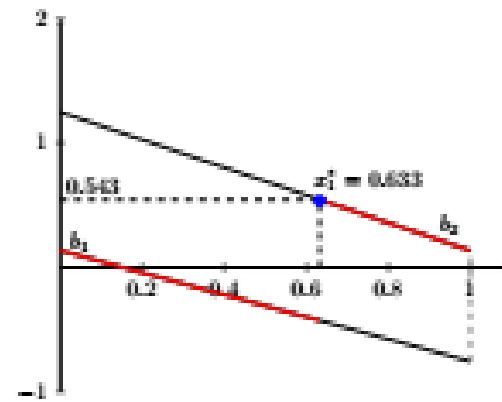


Leader

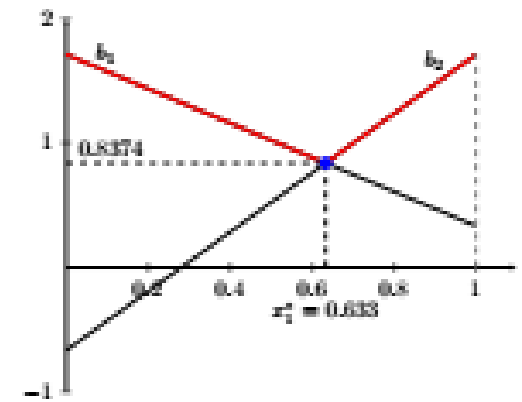


Follower

State s_2



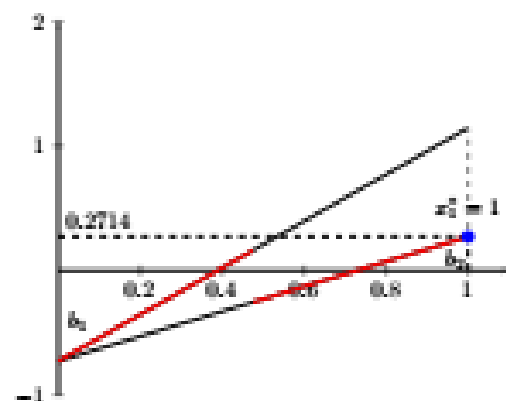
Leader



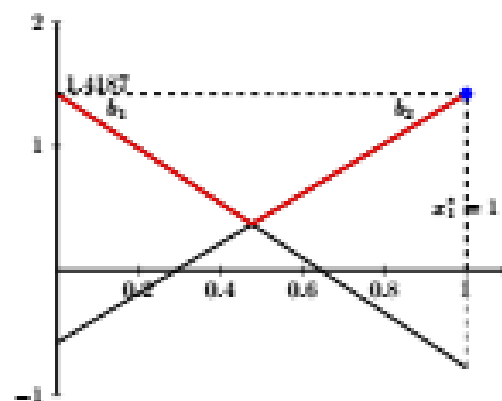
Follower

Iteration 15

State s_1

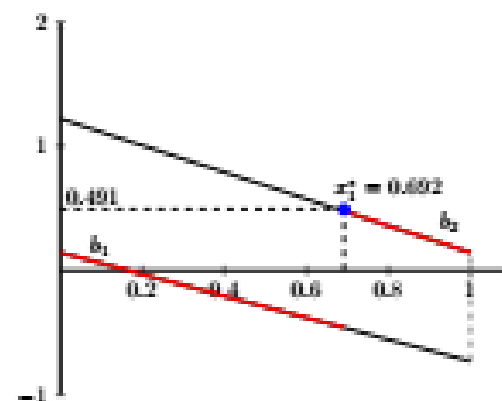


Leader

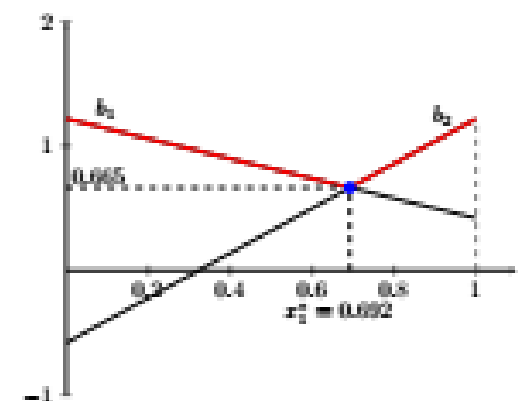


Follower

State s_2



Leader



Follower



Algorithm 4 VI modified: Infinite horizon for the general case

- 1: Initialize with $n = 0$, $v_A^0(s) = v_B^0(s) = 0$ for every $s \in \mathcal{S}$.
- 2: **for** $n = 1, \dots, MAX_IT$ **do**
- 3: Find the pair (v_A^n, v_B^n) by

$$(v_A^n, v_B^n)(s) = T(v_A^{n-1}, v_B^{n-1})(s) .$$

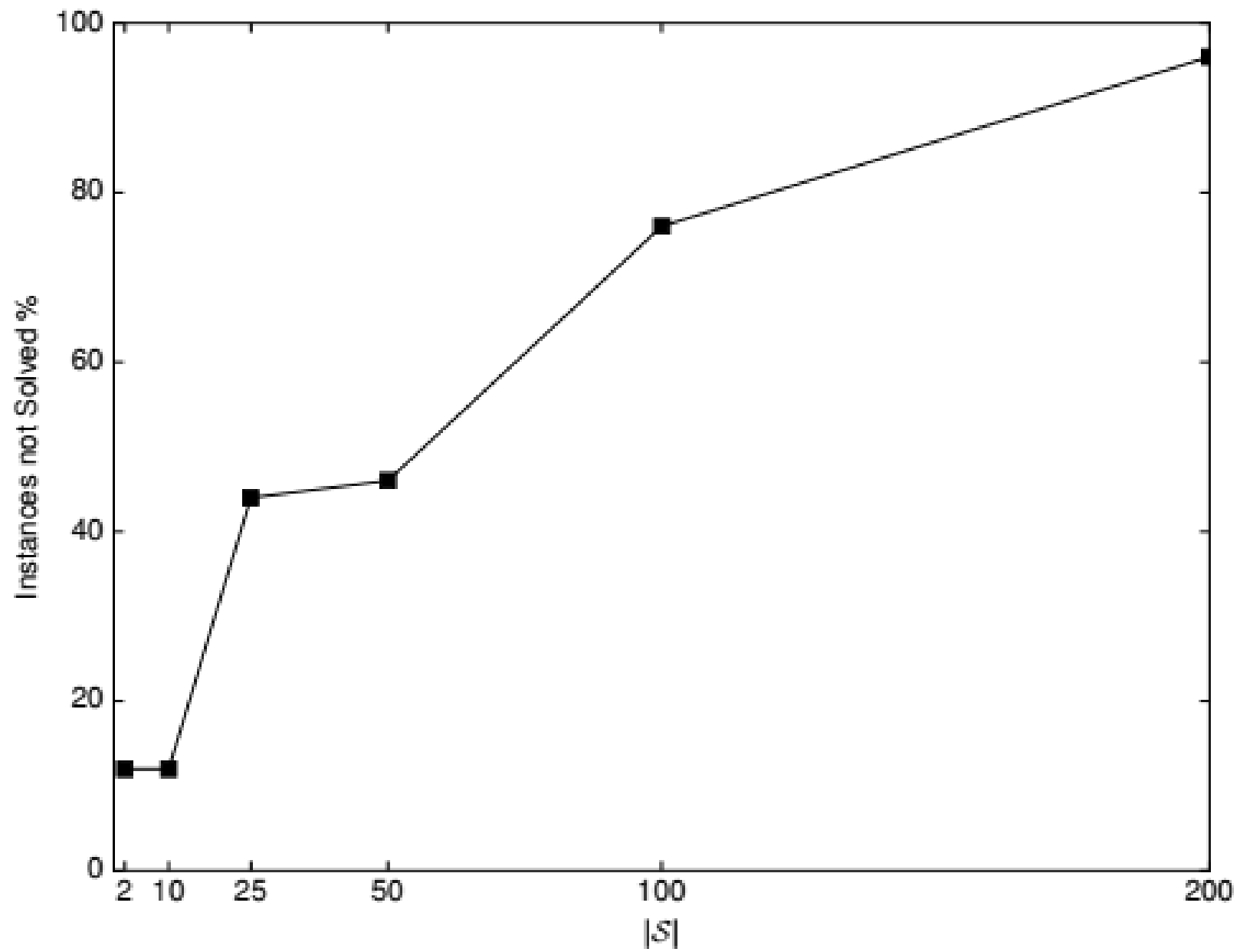
Finding f^* and g^* SSE strategies at stage $n - 1$.

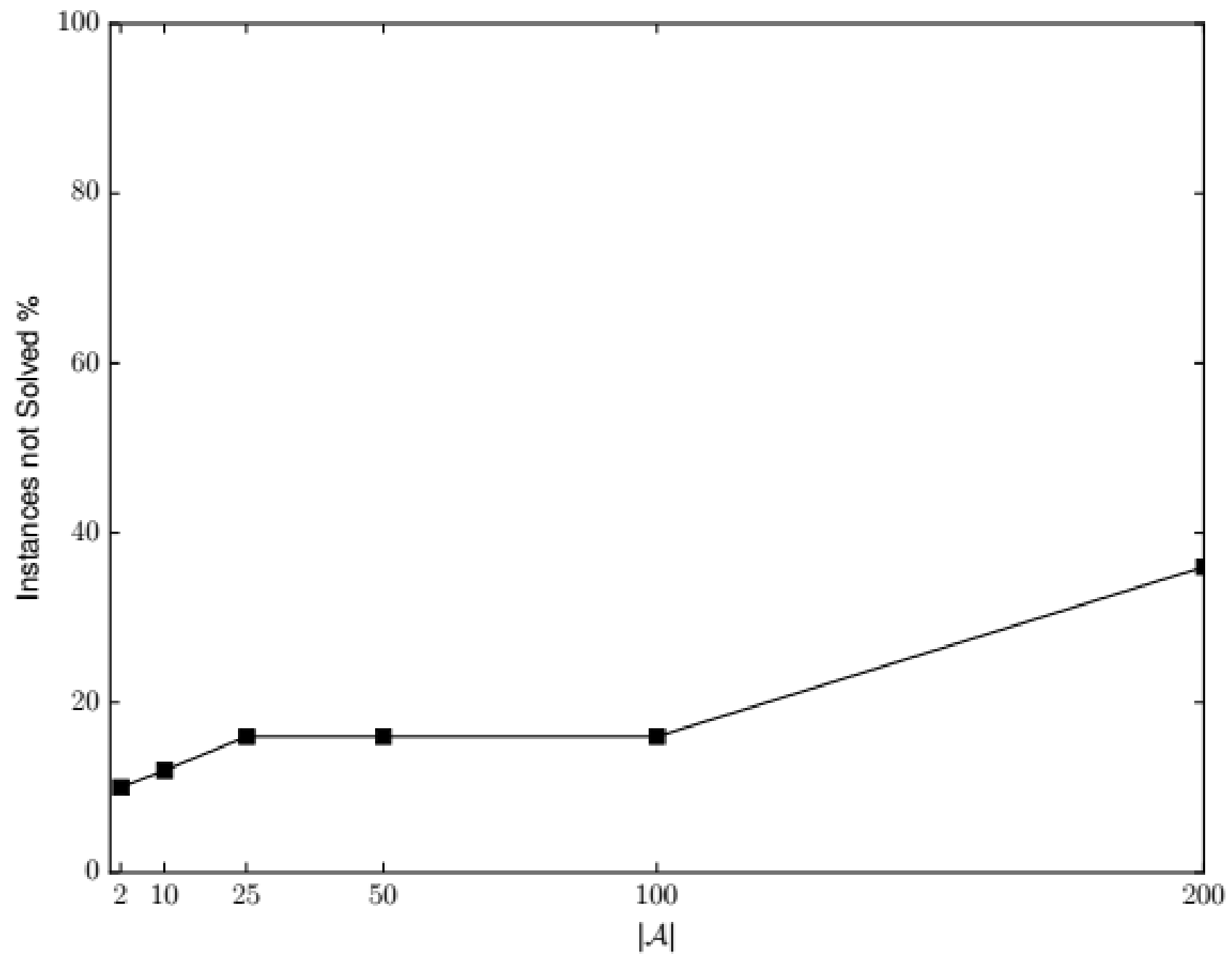
- 4: **if** $(v_A^n, v_B^n) = (v_A^{n-1}, v_B^{n-1})$ **then**
- 5: **return** (v_A^n, v_B^n) fixed point of T .
- 6: **end if**
- 7: **if** $\|(v_A^n, v_B^n) - (v_A^{n-1}, v_B^{n-1})\| > 2 \frac{\beta^{n-1}}{1-\beta} \|(r_A, r_B)\|$ **then**
- 8: **return** UNDEFINED 1.
- 9: **end if**
- 10: **end for**
- 11: **return** UNDEFINED 2.

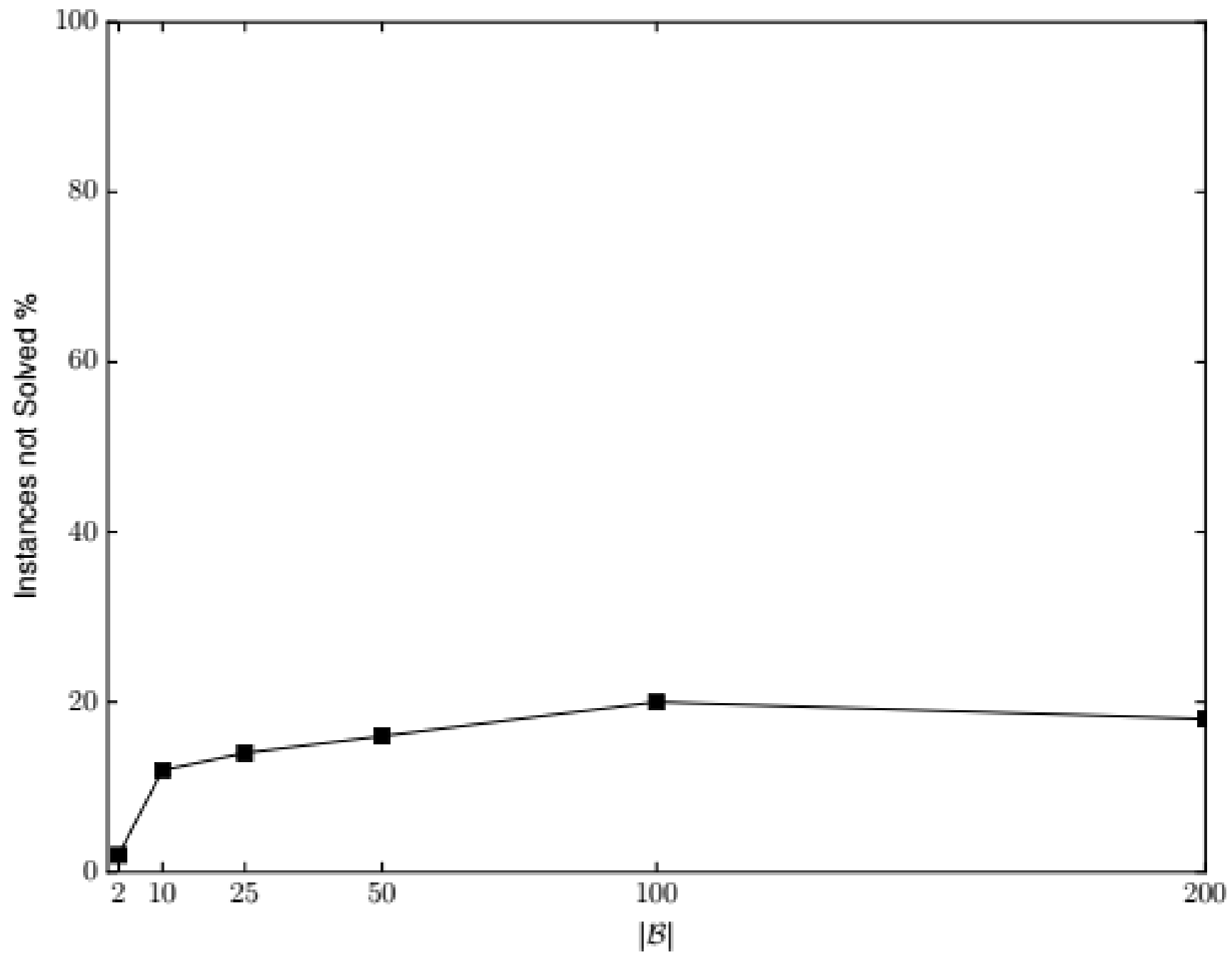
We found an equilibrium.

There is no geometrical decreasing.

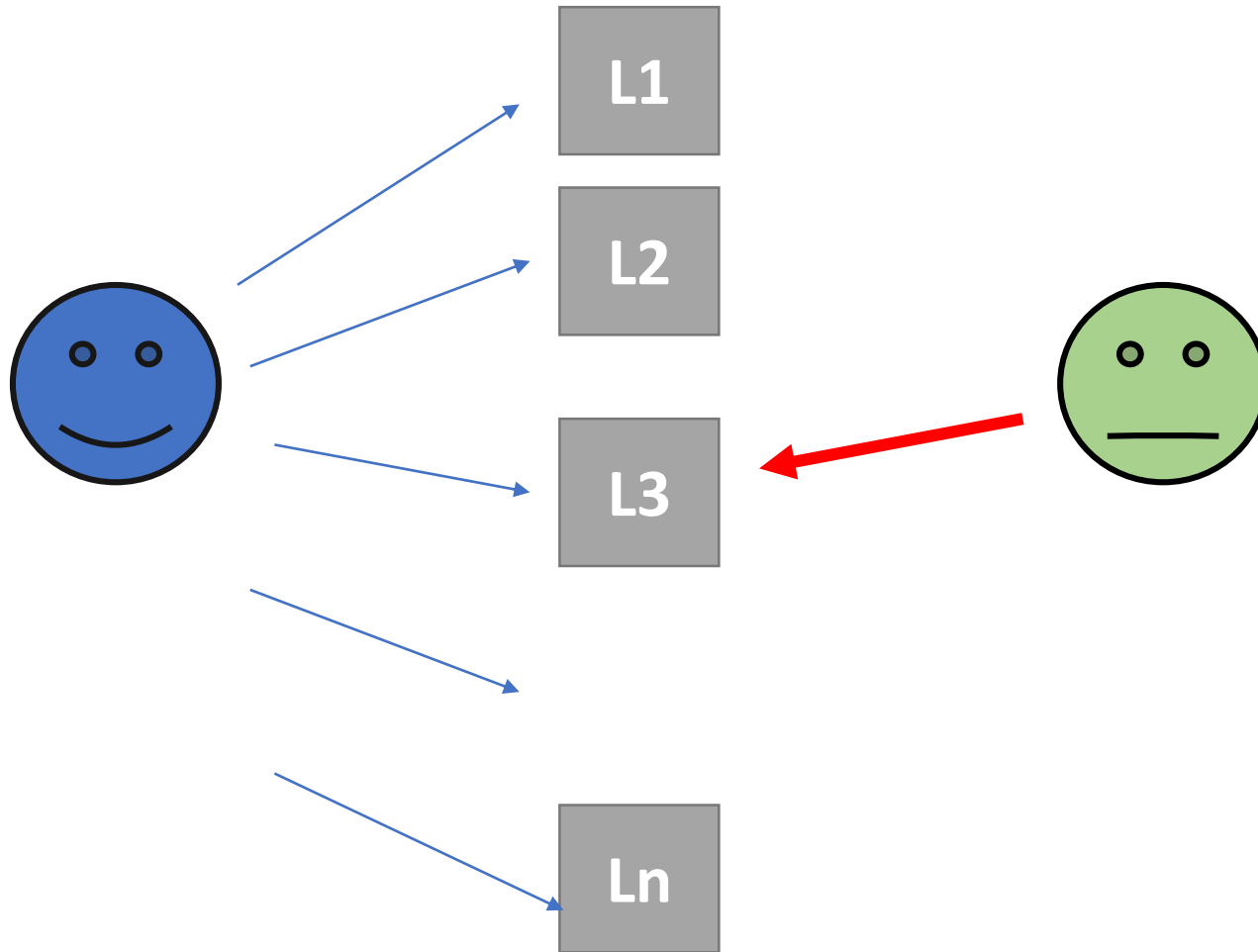
We get tired of finding an equilibrium.







Security Games

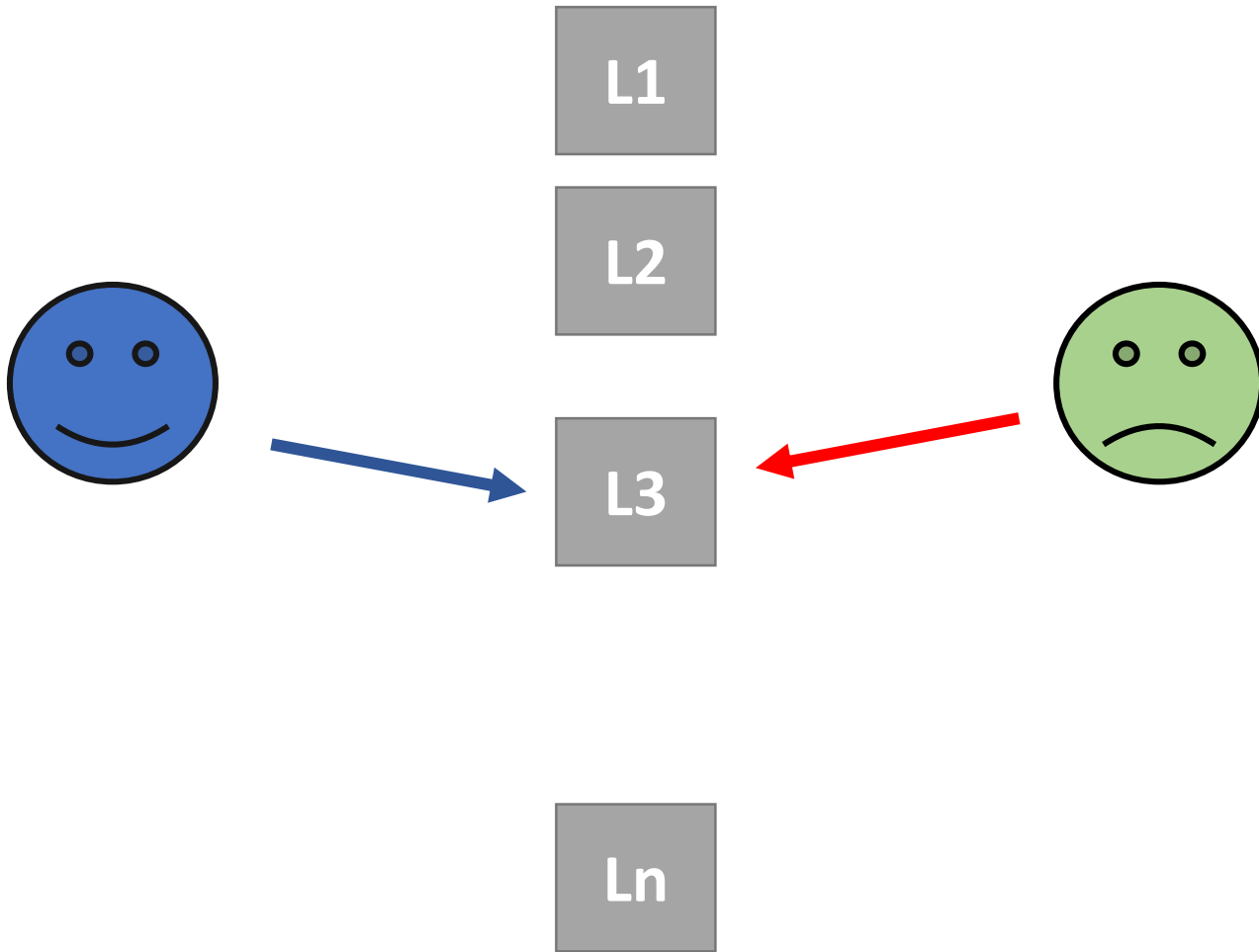


Leader = DEFENDER

Follower = ATTACKER

Payoffs only depends on whether a location is protected or not.

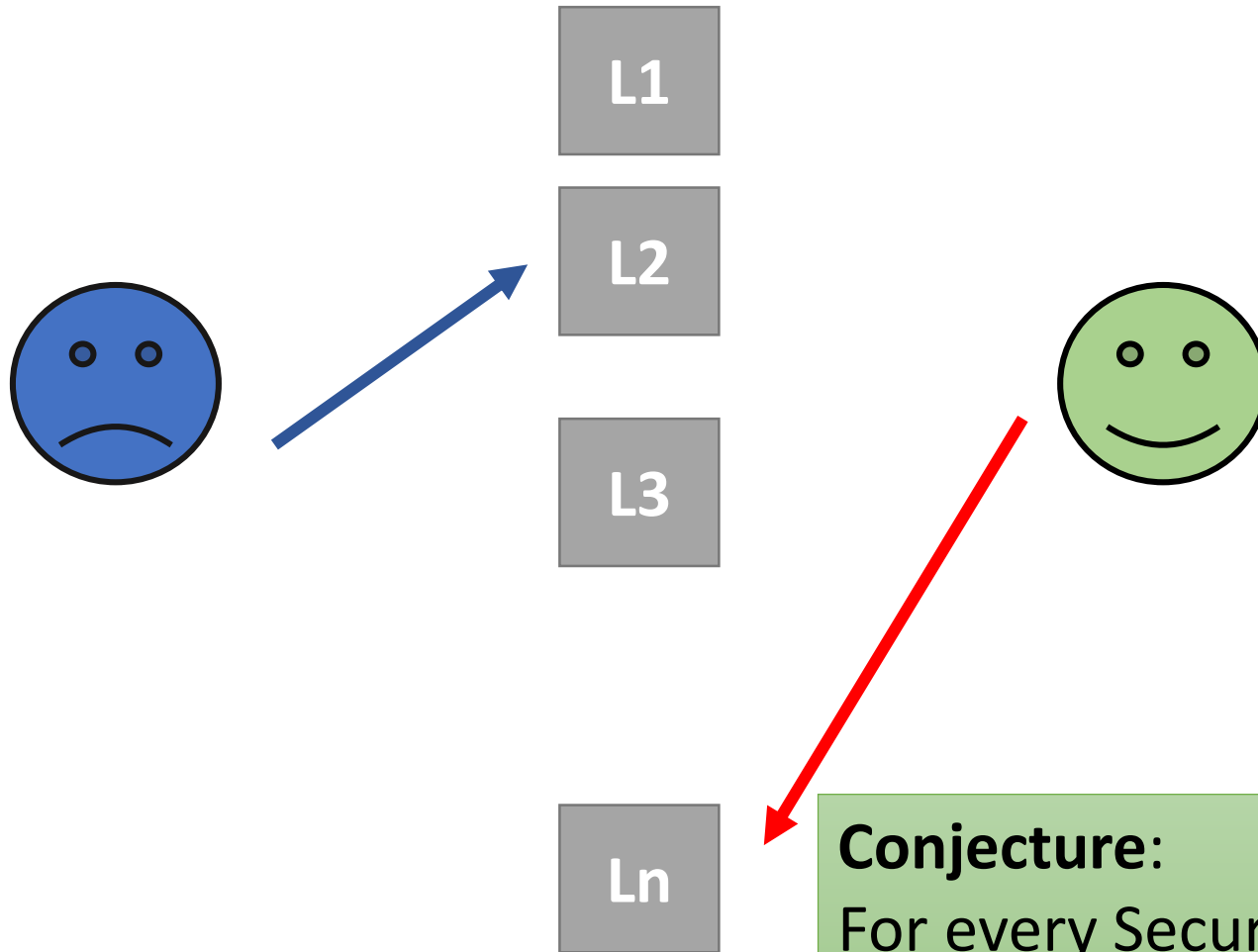
Security Games



Reward_D > 0

Penalty_A < 0

Security Games



Penalty_D < 0

Reward_A > 0

Non pure strategies seems to be optimal for the leader.

Computationally all instances in Security games VI converges with the geometric bound.

Conjecture:

For every Security game with this payoff structure, the operator T is contractive of modulus $\beta = \max\{\beta_A, \beta_B\}$

Conclusions

- We define suitable dynamic programming operators and we use it to prove unicity of values of Strong Stackelberg games in stationary policies for a family of problems.
- We define Value Iteration and Policy Iteration algorithms for finding Stackelberg stationary equilibrium.
- We prove via counterexample that this methodology is not always applicable for the general case.
- We study security games and we conjecture that operators for this type of games are contractive.

Thank You!

Víctor Bucarey López
vbucarey@ing.uchile.cl

INRIA – LILLE
December 2017

