



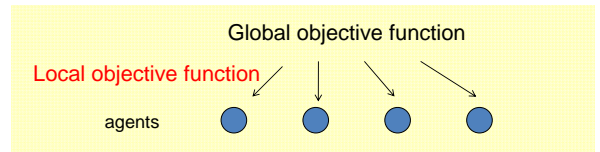
Application of Potential Game for Power Control in Wireless Networks and Network Formation



Goto Tatsuhiko
FL10_04_01
19th, April, 2010



Review (Game theoretic approach)



A_i : Action set $A = \prod_{p_i \in p} A_i$: set of joint action

$a_{-i} = (a_1, a_2, \dots, a_{i-1}, a_{i+1}, \dots, a_n)$
joint action $a = (a_i, a_{-i})$

Local objective function $U_i : A \rightarrow R$

- Control design
1. Designing the player objective function
 2. Learning dynamics (repeated game) (ex)single stage memory dynamics



Review (Potential game)

Global planner $\phi : A \rightarrow \Re$ (potential function)



Make player's objective function U_i

$$U_i(a_i'', a_{-i}) - U_i(a_i', a_{-i}) = \phi_i(a_i'', a_{-i}) - \phi_i(a_i', a_{-i})$$

Changing in the player's objective function = Changing in the potential function

Every agent select an action to maximize their objective function



Outline

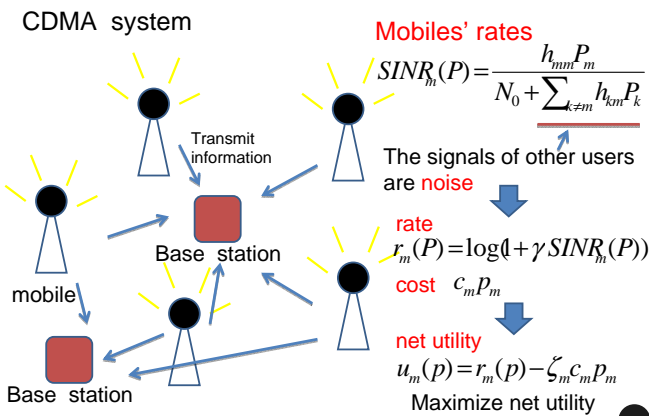
Near-Optimal Power Control in Wireless Networks: A Potential Game Approach

Utku Ozan Candogan, Ihsai Menache, Asuman Ozdaglar and Pablo A. Parrilo. Laboratory for Information and Decision Systems, Massachusetts Institute of Technology, Cambridge, MA, 02139

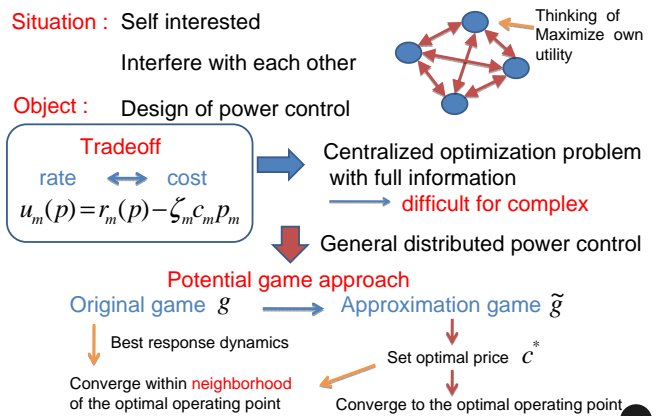
- Background
- Model
- Modified utilities
- Near optimal dynamics
- Convergence analysis
- Simulation Result



background



approach





Outline

Tokyo Institute of Technology

- Background
- **Model**
- Modified utilities
- Near optimal dynamics
- Convergence analysis
- Simulation Result

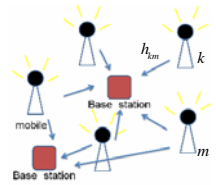
Tokyo Institute of Technology

Fujita Laboratory 7



model

Tokyo Institute of Technology



mobiles **Power allocation**

$$M = \{1, \dots, M\} \quad P = \{p_1, \dots, p_M\}$$

SINR

$$SINR_m(P) = \frac{h_{mm} p_m}{N_0 + \sum_{k \neq m} h_{km} p_k}$$

h_{km} : gain between user k and transmitter m's base station

rate

$$r_m(P) = \log(1 + \gamma SINR_m(P))$$

cost $c_m p_m$

Objective function (net utility)

$$u_m(p) = r_m(p) - \zeta_m c_m p_m$$

p_m : **transmission power**

$$0 \leq p_{\min} \leq p_m \leq \bar{p}_m$$

User specific rate vs money
Tradeoff coefficient

Tokyo Institute of Technology

Fujita Laboratory 8



Power game (definition)

Tokyo Institute of Technology

Power game

$$g = \langle M, \{u_m\}_{m \in M}, \{p_m\}_{m \in M} \rangle$$

m's objective function m's action set

Self interested

$$\max_{\tilde{p}_m \in P_m} u_m(\tilde{p}_m, p_{-m})$$

Nash equilibrium (NE)

$$u_m(p) \geq u_m(\tilde{p}_m, p_{-m}) \quad \forall \tilde{p}_m \in p_m, \forall m \in M$$

ϵ -**Nash equilibrium**

$$u_m(p) \geq u_m(q_m, p_{-m}) - \epsilon \quad \forall q_m \in p_m, \forall m \in M$$

Central planner wishes to impose some **performance objective**

System Utility

$$\max_{p \in P} U_0(p)$$

$$(ex) U_0(p) = \sum_m r_m(p)$$

Sum rate objective

→ Optimal solution p^* (desired operating point)

Tokyo Institute of Technology

Fujita Laboratory 9



Modified utilities

Tokyo Institute of Technology

modified Utility

$$\tilde{u}_m(p) = \tilde{r}_m(p) - \zeta_m c_m p_m$$

$$\tilde{r}_m(p) = \log(\gamma SINR_m(p))$$

$$r_m(P) = \log(1 + \gamma SINR_m(P))$$

$$SINR_m(P) = \frac{h_{mm} p_m}{N_0 + \sum_{k \neq m} h_{km} p_k}$$

Approximation is good
Spreading gain $\gamma \gg 1$
or
 $h_{mm} \gg h_{mk}$

Can make potential function

$$\phi(p) = \sum_m \log(p_m) - \zeta_m c_m p_m$$

$$(\phi(p_m, p_{-m}) - \phi(q_m, p_{-m})) = \tilde{u}_m(p_m, p_{-m}) - \tilde{u}_m(q_m, p_{-m})$$

Strictly concave → **unique NE**

Potential game

$$\tilde{g} = \langle M, \{\tilde{u}_m\}_{m \in M}, \{p_m\}_{m \in M} \rangle$$

Tokyo Institute of Technology

Fujita Laboratory 10



Assigning prices

Tokyo Institute of Technology

Assigning prices c^* to coincide with NE of \tilde{g} and p^*

$$\tilde{u}_m(p) = \tilde{r}_m(p) - \zeta_m c_m p_m$$

[Theorem]

Let p^* be the desired operating point. Then the prices c^* are given by

$$c_m^* = (\zeta_m p_m^*)^{-1} \quad m \in M$$

(proof)

$\phi(p)$ **Strictly concave** → **unique NE**

→ Maxima of $\phi(p)$ is NE

$$\frac{\partial \phi}{\partial p_m} = \frac{1}{p_m} - \frac{1}{p_m^*} \rightarrow p = p^* \rightarrow \frac{\partial \phi}{\partial p_m} = 0$$

$$\phi(p) = \sum_m \log(p_m) - \zeta_m c_m^* p_m$$

$$c_m^* = (\zeta_m p_m^*)^{-1}$$

→ p^* Global maximum of the potential

Tokyo Institute of Technology

Fujita Laboratory 11



Near optimal dynamics

Tokyo Institute of Technology

p^* is not NE of the game g with c^*

→ Converge neighbor of p^* ?

Best Response dynamics

$$p_m \leftarrow p_m + \alpha (\beta_m(p_{-m}) - p_m)$$

$$\text{Best Response } \beta_m(p_{-m}) = \arg \max_{p_m \in P_m} u_m(p_m, p_{-m})$$

α is small

most good action for user m

$$\dot{p}_m = \beta_m(p_{-m}) - p_m$$

\tilde{g} with $c = c^*$ → **Converge to p^***
BR (Lyapunov analysis)

How about g ?

Tokyo Institute of Technology

Fujita Laboratory 12



Outline

Tokyo Institute of Technology

- Background
- Model
- Modified utilities
- Near optimal dynamics
- **Convergence analysis**
- Simulation Result

Tokyo Institute of Technology

Fujita Laboratory 13



Convergence analysis

Tokyo Institute of Technology

$$\begin{aligned} \text{Best Response of } \tilde{g} \quad \tilde{\beta}_m(p_{-m}) &= \arg \max_{p_m \in P_m} \tilde{u}_m(p_m, p_{-m}) \\ &= \arg \max_{p_m \in P_m} \phi(p_m, p_{-m}) \quad (\text{From PG}) \end{aligned}$$

$$\varepsilon\text{-equilibria of } \tilde{g} \quad \tilde{I}_\varepsilon = \{p \mid \tilde{u}_m(p_m, p_{-m}) \geq \tilde{u}_m(q_m, p_{-m}) - \varepsilon\}$$

[Lemma]

$$\text{The BR in } \tilde{g} \text{ converge to } \tilde{I}_\varepsilon \quad \varepsilon \leq \frac{1}{\gamma} \sum_{m \in M} \frac{1}{\text{SINR}_{\min m}}$$

$$\text{SINR}_{\min m}(P) = \frac{h_{mm} P_{\min m}}{N_0 + \sum_{k \neq m} h_{km} P_{\max k}}$$

(proof) $\bar{\phi}$: maximum value of ϕ

$$V = \bar{\phi} - \phi \geq 0 : \text{Lyapunov function}$$

Tokyo Institute of Technology

Fujita Laboratory 14



proof

Tokyo Institute of Technology

$$\begin{aligned} -\dot{V} &= \sum_{m \in M} \frac{\partial \phi}{\partial p_m} (\tilde{\beta}_m(p_{-m}) - p_m) + \sum_{m \in M} \frac{\partial \phi}{\partial p_m} (\beta_m(p_{-m}) - \tilde{\beta}_m(p_{-m})) \\ &\quad \sum_{m \in M} \frac{\partial \phi}{\partial p_m} (\tilde{\beta}_m(p_{-m}) - p_m) \geq \tilde{u}_m(\tilde{\beta}_m(p_{-m}), p_{-m}) - \tilde{u}_m(p_m, p_{-m}) \\ &\quad \left| \frac{\partial \phi}{\partial p_m} (\beta_m(p_{-m}) - \tilde{\beta}_m(p_{-m})) \right| \leq \left(\frac{1}{p_{\min m}} - \frac{1}{\bar{p}_m} \right) \xi_m \\ &\quad \left(\xi_m = \frac{N_0}{h_{mm}} + \sum_{k \neq m} \frac{h_{km}}{h_{mm}} \bar{p}_k \right) \\ -\dot{V} &\geq \sum_{m \in M} (\tilde{u}_m(\tilde{\beta}_m(p_{-m}), p_{-m}) - \tilde{u}_m(p_m, p_{-m})) - \sum_{m \in M} \left(\frac{1}{p_{\min m}} - \frac{1}{\bar{p}_m} \right) \xi_m \\ &\quad \sum_{m \in M} (\tilde{u}_m(\tilde{\beta}_m(p_{-m}), p_{-m}) - \tilde{u}_m(p_m, p_{-m})) \geq \sum_{m \in M} \frac{\xi_m}{\mathcal{P}_{\min m}} \\ &\quad \rightarrow \dot{V} \leq - \sum_{m \in M} \frac{\xi_m}{\mathcal{P}_m} \end{aligned}$$

Tokyo Institute of Technology

Fujita Laboratory 15



proof

Tokyo Institute of Technology

$$\sum_{m \in M} (\tilde{u}_m(\tilde{\beta}_m(p_{-m}), p_{-m}) - \tilde{u}_m(p_m, p_{-m})) \geq \sum_{m \in M} \frac{\xi_m}{\mathcal{P}_{\min m}}$$

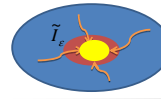
$$\rightarrow \dot{V} \leq - \sum_{m \in M} \frac{\xi_m}{\mathcal{P}_m}$$

Converge to this set from Lyapunov method

$$\tilde{u}_m(\tilde{\beta}_m(p_{-m}), p_{-m}) - \tilde{u}_m(p_m, p_{-m}) \leq \sum_{m \in M} \frac{\xi_m}{\mathcal{P}_{\min m}}$$

$$\rightarrow \varepsilon \leq \frac{1}{\gamma} \sum_{m \in M} \frac{1}{\text{SINR}_{\min m}}$$

$$(\tilde{I}_\varepsilon = \{p \mid \tilde{u}_m(p_m, p_{-m}) \geq \tilde{u}_m(q_m, p_{-m}) - \varepsilon\})$$



$$\dot{V} \leq - \sum_{m \in M} \frac{\xi_m}{\mathcal{P}_m}$$

Tokyo Institute of Technology

Fujita Laboratory 16



How far ?

Tokyo Institute of Technology

How far the set of ε - equilibria of \tilde{g} from p^* ?

[Theorem] $\left| \tilde{p}_m - p_m^* \right| \leq \bar{P}_m \sqrt{2\varepsilon} \quad \tilde{p} \in \tilde{I}_\varepsilon$

ε -equilibria of \tilde{g} \rightarrow Global maximum of the potential

$$\tilde{I}_\varepsilon = \{p \mid \tilde{u}_m(p_m, p_{-m}) \geq \tilde{u}_m(q_m, p_{-m}) - \varepsilon\}$$

(proof)

$$\phi(p_m^*, \tilde{p}_{-m}) - \phi(\tilde{p}_m, \tilde{p}_{-m}) \leq \varepsilon$$

$$\rightarrow (\log(p_m^*) - \lambda_m p_m^*) - (\log(\tilde{p}_m) - \lambda_m \tilde{p}_m) \leq \varepsilon \quad \left(\phi(p) = \sum_m \log(p_m) - \sum_m \lambda_m p_m \right)$$

$$\rightarrow f_m(p_m^*) - f_m(\tilde{p}_m) \leq \varepsilon \quad \left(f_m = \log(p_m) - \lambda_m p_m \right)$$

Tokyo Institute of Technology

Fujita Laboratory 17



proof

Tokyo Institute of Technology

$$f_m = \log(p_m) - \lambda_m p_m$$

$$\rightarrow f_m(\tilde{p}_m) = f_m(p_m^*) + (\tilde{p}_m - p_m^*) \frac{\partial f_m(p_m^*)}{\partial p_m} + \frac{1}{2} (\tilde{p}_m - p_m^*)^2 \frac{\partial^2 f_m(p_m^* + \alpha(\tilde{p}_m - p_m^*))}{\partial p_m^2}$$

p^* Is desired operating point

$$\rightarrow \frac{\partial f_m(p_m^*)}{\partial p_m} = \frac{\partial \phi(p_m^*)}{\partial p_m} = 0$$

$$\rightarrow f_m(p_m^*) - f_m(\tilde{p}_m) = \frac{1}{2} (p_m^* - \tilde{p}_m)^2 \frac{1}{(p_m^* + \alpha(\tilde{p}_m - p_m^*))^2}$$

$$\rightarrow 2(p_m^* + \alpha(\tilde{p}_m - p_m^*))^2 (f_m(p_m^*) - f_m(\tilde{p}_m)) = (p_m^* - \tilde{p}_m)^2$$

$$\rightarrow 2\varepsilon \bar{P}_m^2 \geq (p_m^* - \tilde{p}_m)^2 \quad (f_m(p_m^*) - f_m(\tilde{p}_m) \leq \varepsilon, 0 < p_m^*, \tilde{p}_m \leq \bar{p}_m)$$

$$\rightarrow \left| \tilde{p}_m - p_m^* \right| \leq \bar{P}_m \sqrt{2\varepsilon}$$

Tokyo Institute of Technology

Fujita Laboratory 18



Near optimal performance

Tokyo Institute of Technology

Near optimal performance in terms of **system utility**

Performance loss decrease with small ε
increase with large L, L_m

$$(ex) U_0(p) = \sum_m r_m(p)$$

Sum rate objective

[Theorem] Let $\varepsilon \leq \frac{1}{\gamma} \sum_{m \in M} \frac{1}{SINR_{\min m}}$ then

(1) U_0 is Lipschitz continuous function, with L . Then

$$|U_0(p^*) - U_0(\tilde{p})| \leq \sqrt{2\varepsilon} L \sqrt{\sum_{m \in M} \bar{P}_m^2}$$

(2) Assume that U_0 is a continuous differentiable function such that

$$\left| \frac{\partial U_0}{\partial p_m} \right| \leq L_m \quad \text{Then} \quad |U_0(p^*) - U_0(\tilde{p})| \leq \sqrt{2\varepsilon} \sum_{m \in M} \bar{P}_m L_m$$

Difference between p^* and \tilde{p} \rightarrow Difference between $U_0(p^*)$ and $U_0(\tilde{p})$

Tokyo Institute of Technology

Fujita Laboratory 19



Outline

Tokyo Institute of Technology

- Background
- Model
- Modified utilities
- Near optimal dynamics
- Convergence analysis
- **Simulation Result**

Tokyo Institute of Technology

Fujita Laboratory 20



Simulation result

Tokyo Institute of Technology

• Three users

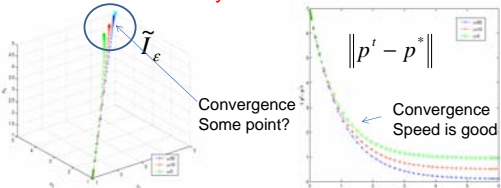
• desired operating point $p^* = [5, 5, 5]$

$r_m(p) = \log(1 + \gamma SINR_m(p))$ $\gamma \in \{5, 10, 15\}$

$$SINR_m(p) = \frac{h_{mm} P_m}{N_0 + \sum_{k \neq m} h_{km} P_k}$$

$N_0 = 1$ $h_{km} \in [0, 2]$ $h_{mm} \in [2, 4]$

Do the BR dynamics with c^*



\rightarrow Big γ is good

$$\varepsilon \leq \frac{1}{\gamma} \sum_{m \in M} \frac{1}{SINR_{\min m}}$$

\rightarrow Monotonously decreasing

Tokyo Institute of Technology

Fujita Laboratory 21



Simulation result (system utility)

Tokyo Institute of Technology

Sum rate objective

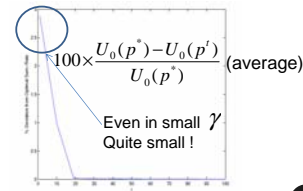
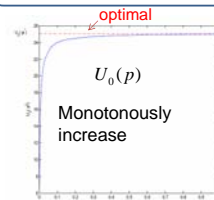
$$U_0(p) = \sum_m r_m(p) \quad \rightarrow \quad \max_{p \in P} U_0(p)$$

[Theorem]

U_0 is a continuous differentiable function

$$\left| \frac{\partial U_0}{\partial p_m} \right| \leq \frac{M-1}{p_{\min m}}$$

$$\rightarrow |U_0(p^*) - U_0(\tilde{p})| \leq \sqrt{2\varepsilon} (M-1) \sum_{m \in M} \frac{\bar{P}_m}{P_{\min m}}$$
$$\left(|U_0(p^*) - U_0(\tilde{p})| \leq \sqrt{2\varepsilon} \sum_{m \in M} \bar{P}_m L_m \right)$$



Tokyo Institute of Technology

Fujita Laboratory 22



General type

Tokyo Institute of Technology

$$u_m \rightarrow \phi \rightarrow \hat{u}_m$$

Find the most close potential function of the game by u_m

Find the most close objective function of ϕ to u_m

$$d^2(g) = \min_{\phi \in C_0} \left\| \delta_0 \phi - \sum_{m \in M} D_m u^m \right\|_2^2$$

$$\hat{u}^m = \operatorname{argmin}_{\bar{u}^m} \left\| u^m - \bar{u}^m \right\|_2^2$$

$$D_m \bar{u}_m = D_m \phi$$

D_m : Difference

$$W^n(p, q) = 1 \quad p \neq q$$

$$(D_m \phi)(p, q) = W^n(p, q)(\phi(q) - \phi(p)) \quad W^n(p, q) = 0 \quad p = q$$

Potential game

$$D_m u_m = D_m \phi$$

$$\sum_{m \in M} D_m u_m = \sum_{m \in M} D_m \phi = \delta_0 \phi$$

combinatorial gradient operator $\delta_0 = \sum_{m \in M} D_m$

Tokyo Institute of Technology

Fujita Laboratory 23



Theorem of Projection

Tokyo Institute of Technology

[Theorem] Optimal projection

$$\phi = \left(\sum_{m \in M} \Pi_m \right)^\dagger \sum_{m \in M} \Pi_m u^m$$

$$\hat{u}^m = (I - \Pi_m) u^m + \Pi_m \left(\sum_{k \in M} \Pi_k \right)^\dagger \sum_{k \in M} \Pi_k u^k \quad \Pi_m = D_m^* D_m$$

[Theorem]

Any equilibrium of \tilde{g} is an ε -equilibrium of g .

$$\varepsilon \leq \sqrt{2} d(g) \quad d^2(g) = \min_{\phi \in C_0} \left\| \delta_0 \phi - \sum_{m \in M} D_m u^m \right\|_2^2$$

g : game

\tilde{g} : projection of the game g

$$u_m(p) \geq u_m(q_m, p_{-m}) - \varepsilon$$

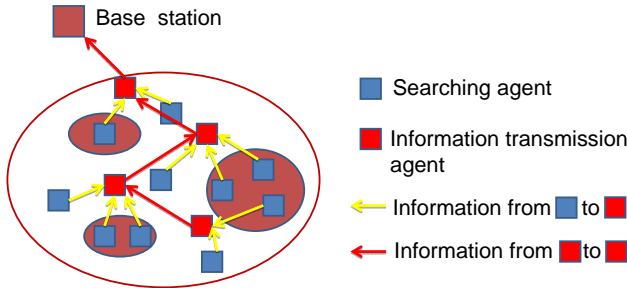
Tokyo Institute of Technology

Fujita Laboratory 24



Coverage considering wireless network

Tokyo Institute of Technology



Object : **optimal coverage** and make **small** transmission cost of information
 → Several Information transmission agent needed
 Agent can change from **blue** to **red** and from **red** to **blue**

Tokyo Institute of Technology

Fujita Laboratory 25



Outline

Tokyo Institute of Technology

Distributed Dynamic Reinforcement of Efficient Outcomes in Multiagent Coordination and Network Formation*

Georgios C. Chasparis¹ and Jeff S. Shamma¹

November 7, 2009

- **Background**
- Reinforcement learning
- Asymptotic stability analysis
- Dynamic Reinforcement
- Asymptotic stability of RADR
- Simulation Result

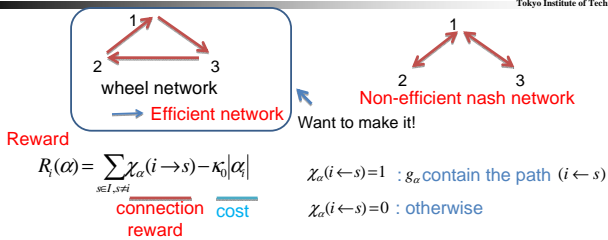
Tokyo Institute of Technology

Fujita Laboratory 26



Distributed network formation

Tokyo Institute of Technology



Reward $R_i(\alpha) = \sum_{s \in I, s \neq i} \chi_{\alpha}(i \rightarrow s) - k_0 |\alpha_i|$ $\chi_{\alpha}(i \leftarrow s) = 1$: g_{α} contain the path $(i \leftarrow s)$
 connection cost reward $\chi_{\alpha}(i \leftarrow s) = 0$: otherwise

Actions of agent
 (ex) $A_i = \{\emptyset, \{2\}, \{3\}, \{2,3\}\}$

Nash network α^* Nash network $\iff R_i(\alpha_i^*, \alpha_{-i}^*) \geq R_i(\alpha_i, \alpha_{-i}^*) \quad \alpha_i^* \in A_i \setminus \alpha_i^*$
 wheel network
 Connected network is uniquely defined by a path $(i \leftarrow i)$
 → Every agent realizes its maximum possible utility

Tokyo Institute of Technology

Fujita Laboratory 27

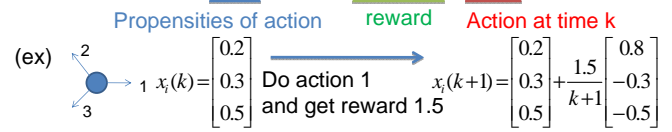


Reinforcement learning

Tokyo Institute of Technology

Learning algorithm

$$x_i(k+1) = x_i(k) + \varepsilon(k) R_i(\alpha(k)) [\alpha_i(k) - x_i(k)]$$



Probability of action selection

$$\sigma_i(k) = (1-\lambda)x_i(k) + \lambda \frac{1}{|A_i|} \mathbf{1}$$

$\varepsilon(k) = \frac{1}{k+1}$
 Propensities are proportional to the cumulative reward

Random selection
 $x_i(k)$ evolves over the probability simplex $\Delta(A_i)$ ($\Delta(m) = \{v \in \mathfrak{R}^m \mid v \geq 0, \mathbf{1}^T v = 1\}$)

Tokyo Institute of Technology

Fujita Laboratory 28



Asymptotic stability analysis

Tokyo Institute of Technology

$$x_i(k+1) = x_i(k) + \varepsilon(k) R_i(\alpha(k)) [\alpha_i(k) - x_i(k)]$$

$$\bar{g}_i(x(k)) = \bar{r}_i(x(k)) - \bar{R}_i(x(k)) x_i(k)$$

$$\bar{r}_i(x(k)) = E[R_i(\alpha(k)) \alpha_i(k) \mid x(k)] \quad \bar{R}_i(x(k)) = E[R_i(\alpha(k)) \mid x(k)]$$

$$\xi_i(k) = R_i(\alpha(k)) [\alpha_i(k) - x_i(k)] - \bar{g}_i(x(k))$$

$$x_i(k+1) = x_i(k) + \varepsilon(k) [\bar{g}_i(x(k)) + \xi_i(k)]$$

[proposition3.1] ($\lambda > 0$) $\xrightarrow{\text{deterministic noise}}$
 $x(k+1) = x(k) + \varepsilon(k) [\bar{g}(x(k)) + \xi(k)] \quad x(k) = \{x_1(k), \dots, x_n(k)\}$

$x(k)$ convergence to an **invariant set** of the $\dot{x} = \bar{g}(x)$
 $A \subset \Delta$ be a **locally Asymptotically stable set** → $\Pr[\lim_{k \rightarrow \infty} x(k) \in A] > 0$

Tokyo Institute of Technology

Fujita Laboratory 29



Asymptotic stability analysis

Tokyo Institute of Technology

$\dot{x} = \bar{g}(x)$ → **Stationary point**
 $S = \{x \in X \mid \bar{g}(x) = \bar{r}(x) - \bar{R}(x)x = 0\}$

[proposition3.2] ($\lambda > 0$)

x^* linearly unstable stationary point of $\dot{x} = \bar{g}(x)$

$$\Pr[\lim_{k \rightarrow \infty} x(k) = x^*] = 0$$

Expected reward $\bar{v}_i(j, x^*) = E[R_i(\alpha) \mid \alpha_i = j, x_{-i} = x_{-i}^*] \quad x^* \in S$

[proposition3.3] (Stationary point)

x^* stationary point of $\dot{x} = \bar{g}(x)$ ↔ $\bar{v}_i(j, x^*) = c_i \quad j \in A_i$
 Expected reward not change by one agent

Tokyo Institute of Technology

Fujita Laboratory 30



proposition

Tokyo Institute of Technology

[proposition3.4] (pure strategies) ($\lambda=0$)

pure strategies profile \rightarrow stationary point of $\dot{x}=\bar{g}(x)$

[proposition3.5] (sensitivity of pure strategies) ($\lambda>0$)

x^* pure strategies profile \rightarrow Exist a differentiable function $v^*(\lambda)$, such $\tilde{x}=x^*+v^*(\lambda)$ and strict NE stationary point of $\dot{x}=\bar{g}(x)$ ($\lim_{\lambda \rightarrow 0} v^*(\lambda)=v^*(0)=0$)

[proposition3.6] (LAS) ($\lambda>0$)

\tilde{x} Locally asymptotically Stable point of $\dot{x}=\bar{g}(x)$ $\leftrightarrow \bar{v}_i(j^*, \tilde{x}) > \bar{v}_i(s, \tilde{x}) \quad \forall s \in A_j \setminus j^*$
 $x_i^* = e_j$

[proposition3.7] ($\lambda>0$)

$\bar{v}_i(j^*, \tilde{x}) > \bar{v}_i(s, \tilde{x}) \quad \forall s \in A_j \setminus j^*$:strict NE $\rightarrow \Pr[\lim_{k \rightarrow \infty} x(k) = \tilde{x}] > 0$ (from p3.6)

$\bar{v}_i(j^*, \tilde{x}) < \bar{v}_i(s, \tilde{x}) \quad \exists s \in A_j \setminus j^*$:not NE $\rightarrow \Pr[\lim_{k \rightarrow \infty} x(k) \in B_\delta(x^*)] = 0$

Neighbor of x^* 31

Tokyo Institute of Technology

Fujita Laboratory



Outline

Tokyo Institute of Technology

- Background
- Reinforcement learning
- Asymptotic stability analysis
- **Dynamic Reinforcement**
- Asymptotic stability of RADR
- Simulation Result

Tokyo Institute of Technology

Fujita Laboratory 32



Dynamic Reinforcement

Tokyo Institute of Technology

Before depend only on the probability distribution \rightarrow Also affected by the history of x_i

\rightarrow Goal is to investigate the effects on convergence to an Efficient pure equilibrium

Learning algorithm

$$x_i(k+1) = x_i(k) + \varepsilon(k) R_i(\alpha(k)) [\alpha_i(k) - x_i(k)]$$

Propensities of action \rightarrow reward \rightarrow Action at time k

Probability of action selection

$$\sigma_i(k) = \prod_{s \in \Delta(m)} [(1-\lambda)(x_i(k) + u_i(k)) + \lambda \frac{1}{|A_i|}]^{\mathbb{1}_{\{s=x\}}}$$

Correspond to history of x_i

$u_i(k) = \gamma_i(\rho_i(k))(x_i(k) - y_i(k))$ $\gamma_i(\rho_i(k))$: RADR parameter

$y_i(k+1) = y_i(k) + \varepsilon(k)(x_i(k) - y_i(k))$ Running average of x_i

$\rho_i(k+1) = \rho_i(k) + \varepsilon(k)(R_i(\alpha(k)) - \rho_i(k))$ Running average of $R_i(\alpha(k))$

Tokyo Institute of Technology

Fujita Laboratory 33



Asymptotic stability of RADR

Tokyo Institute of Technology

$$z(k) = \begin{pmatrix} x(k) \\ y(k) \\ \rho(k) \end{pmatrix} \xrightarrow{\text{Relevant ODE}} \begin{pmatrix} \dot{x} \\ \dot{y} \\ \dot{\rho} \end{pmatrix} = \begin{pmatrix} \bar{g}(z) \\ x-y \\ \bar{R}(z) - \rho \end{pmatrix} \xrightarrow{\text{linearization}} \frac{d}{dt} \begin{pmatrix} \delta x(t) \\ \delta y(t) \\ \delta \rho(t) \end{pmatrix} = \tilde{A} z \begin{pmatrix} \delta x(t) \\ \delta y(t) \\ \delta \rho(t) \end{pmatrix}$$

$$\begin{cases} x_i(t) = \tilde{x}_i + N \delta x_i(t) \\ y_i(t) = \tilde{y}_i + N \delta y_i(t) \\ \delta \rho(t) = \rho(t) - \tilde{\rho} \end{cases}$$

[Theorem 4.1] (LAS of RADR)

RADR parameter

equilibrium $\tilde{z} = (\tilde{x}, \tilde{y}, \tilde{\rho})$ is LAS point of linearization $\leftrightarrow 0 \leq \gamma_i(\tilde{\rho}_i(k)) \leq \frac{\bar{v}_i(j^*, \tilde{x}) + 1}{\bar{v}_i(s, \tilde{x})} - 1 \quad \forall s \neq j^*$

[Theorem 4.2] (RADR convergence)

$0 \leq \gamma_i(\tilde{\rho}_i(k)) \leq \frac{\bar{v}_i(j^*, \tilde{x}) + 1}{\bar{v}_i(s, \tilde{x})} - 1 \quad \forall s \neq j^*, \forall i \rightarrow \Pr[\lim_{k \rightarrow \infty} x(k) = \tilde{x}] > 0$

$\gamma_i(\tilde{\rho}_i(k)) \geq \frac{\bar{v}_i(j^*, \tilde{x}) + 1}{\bar{v}_i(s, \tilde{x})} - 1 \quad \exists s \neq j^*, \exists i \rightarrow \Pr[\lim_{k \rightarrow \infty} x(k) = \tilde{x}] = 0$

Tokyo Institute of Technology

Fujita Laboratory 34



example

Tokyo Institute of Technology

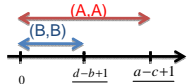
1.A $\begin{matrix} 2A & 2B \\ a,a & b,c \end{matrix}$ (case1) $a > c, d > b, (a-c) > (d-b), a < d$

1.B $\begin{matrix} c,b & d,d \end{matrix}$

Symmetric coordination game

\rightarrow (A,A) risk dominant
(B,B) payoff dominant

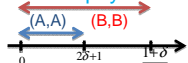
$$0 \leq \gamma_i(\tilde{\rho}_i(k)) \leq \frac{\bar{v}_i(j^*, \tilde{x}) + 1}{\bar{v}_i(s, \tilde{x})} - 1 \quad \forall s \neq j^*$$



$\rightarrow \frac{d-b+1}{b} < \gamma < \frac{a-c+1}{c}$ Positive probability of convergence to (A,A)
Zero probability of convergence to (B,B)

(case2) $a = 2 + 2\delta, b = 1 - \delta, c = 2, d = 1$

\rightarrow (A,A) risk and payoff dominant



$\rightarrow \frac{2\delta+1}{2} < \gamma < \frac{\delta+1}{1-\delta}$ Positive probability of convergence to (B,B)
Zero probability of convergence to (A,A)

Tokyo Institute of Technology

Fujita Laboratory 35



Example 2

Tokyo Institute of Technology

1.A $\begin{matrix} 2A & 2B \\ a_1, a_2 & b, c \end{matrix}$ $a_1 > a_2 > 0, d_1 = a_2, d_2 = a_1, b = c > 0, a_1 > c, d_1 > b$

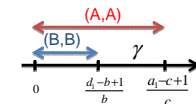
1.B $\begin{matrix} c,b & d_1, d_2 \end{matrix}$

Asymmetric coordination game

\rightarrow (A,A) Desirable for agent 1 ($a_1 < a_2$)
(B,B) Desirable for agent 2 ($a_1 > a_2$)

Only agent 1

$$\frac{d_1 - b + 1}{b} < \gamma < \frac{a_1 - c + 1}{c}$$



Positive probability of convergence to (A,A)
Zero probability of convergence to (B,B)

\rightarrow The agent that applies RADR destabilizes the less Desirable equilibrium in favor of the desirable one

Tokyo Institute of Technology

Fujita Laboratory 36



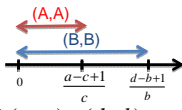
Payoff dependent γ

Tokyo Institute of Technology

	2.A	2.B
1.A	a,a	b,c
1.B	c,b	d,d

Symmetric coordination game

$$a > c > 0, d > b > 0, (a-c) < (d-b), a > d$$



Risk dominant but not
Pay off dominant **can't**
destabilized with constant γ

(A,A) **payoff dominant**
(B,B) **risk dominant**
can't be destabilized

$$\gamma_i(\rho_i) = \frac{\gamma_0}{\rho_i^k} \quad \text{RADR parameter}$$

$$\kappa > \frac{\log\left(\frac{a-c+1}{d-b+1}\right)}{\log\left(\frac{d}{a}\right)} \quad d^k \frac{d-b+1}{b} < \gamma_0 < a^k \frac{a-c+1}{c}$$

Positive probability of convergence to (A,A)

Zero probability of convergence to (B,B)

➔ The risk dominant equilibrium is no longer stable in favor of the payoff dominant equilibrium

Tokyo Institute of Technology

Fujita Laboratory 37



Outline

Tokyo Institute of Technology

- Background
- Reinforcement learning
- Asymptotic stability analysis
- Dynamic Reinforcement
- Asymptotic stability of RADR
- **Simulation Result**

Tokyo Institute of Technology

Fujita Laboratory 38



simulation

Tokyo Institute of Technology

3 agent

Actions of agent

$$A_i = \{A, B, C, D\}$$

$$A_1 = \{\{1\}, \{2\}, \{3\}, \{2,3\}\}$$

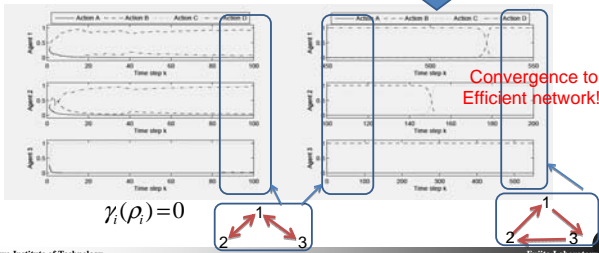
$$A_2 = \{\{1\}, \{2\}, \{3\}, \{1,3\}\}$$

$$A_3 = \{\{1\}, \{2\}, \{3\}, \{2,1\}\}$$

$$\gamma_i(\rho_i) = \frac{\gamma_0}{\rho_i^k} \quad \text{RADR parameter}$$

$$\kappa_0 = \frac{1}{2} \quad \gamma_i \in (2/3, 3/2)$$

➔ Non-efficient Nash network is unstable



Convergence to
Efficient network!

$$\gamma_i(\rho_i) = 0$$

Fujita Laboratory 39

Tokyo Institute of Technology