

非均衡动态博弈：一种智能分层架构

组会报告

谭浚楷

西安交通大学电气学院

2023 年 4 月 22 日



① 研究背景

② 研究内容

③ 计划进度

Systems & Control Letters 125 (2019) 59–66



ELSEVIER

Contents lists available at ScienceDirect

Systems & Control Letters

journal homepage: www.elsevier.com/locate/sysconleNon-equilibrium dynamic games and cyber–physical security: A cognitive hierarchy approach[☆]Aris Kannelopoulos^{*}, Kyriakos G. Vamvoudakis

Daniel Guggenheim School of Aerospace Engineering, Georgia Institute of Technology, Atlanta, GA 30332, USA

- 纳什均衡的背景是绝对理性下的理想情形
- 大多数博弈中，都会有人的存在，此时很可能无法达到纳什均衡
- 同时智能体的计算能力有限，在复杂环境下无法获取无穷智能模型

① 研究背景

② 研究内容

Level- k 模型介绍

零和博弈

认知层分层结构

Model-free Q-learning

③ 计划进度

① 研究背景

② 研究内容

Level- k 模型介绍

零和博弈

认知层分层结构

Model-free Q-learning

③ 计划进度

Level- k 模型介绍

Level- k 模型是什么？

Level- k 智力等级的玩家“自我假设”其他所有对手的智力等级都分布在 Level-0 \sim Level- $(k-1)$ ，且一般服从泊松分布。¹

实际中智力等级应该是怎样的？

实际中应该有无穷个智力等级，其中最理智的模型（即 $k \rightarrow \infty$ ）对应着的最优控制策略；最一般的模型（即 $k=0$ ）对应着的最一般策略（ $u=0$ 或者 $u=random$ ）

¹camerer_cognitive_2004.

① 研究背景

② 研究内容

Level- k 模型介绍

零和博弈

认知层分层结构

Model-free Q-learning

③ 计划进度

Preliminaries

系统：零和博弈（控制 u + 干扰 d ）

$$\dot{x} = Ax + Bu + Kd \quad (1)$$

Value function：价值函数

$$J(x(0), u, d) = \frac{1}{2} \int_0^{\infty} \left(x^T Q x + u^T R u - \gamma^2 \|d\|^2 \right) dt \quad (2)$$

$$V^*(x(t)) = \min_u \max_d \int_t^{\infty} \frac{1}{2} \left(x^T Q x + u^T R u - \gamma^2 \|d\|^2 \right) dt \quad (3)$$

Preliminaries

微分后得到 Hamiltatian

$$H\left(x, \frac{\partial V}{\partial x}, u, d\right) \equiv x^T Qx + u^T Ru - \gamma^2 \|d\|^2 + \left(\frac{\partial V}{\partial x}\right)^T (Ax + Bu + Kd)$$

极值条件

$$0 = \frac{\partial H}{\partial u} \Rightarrow u = -R^{-1} B^T \frac{\partial V}{\partial x} \quad 0 = \frac{\partial H}{\partial d} \Rightarrow u = \frac{1}{\gamma^2} K^T \frac{\partial V}{\partial x} \quad (4)$$

① 研究背景

② 研究内容

Level- k 模型介绍

零和博弈

认知层分层结构

Model-free Q-learning

③ 计划进度

智能等级

Level-0: 初始化选择

Level-0 可以用两种方法模拟：1、认为不存在对手；2、随机行为。本文选择第一种方法

$$V_u^0(x_0) = \min_u \int_0^\infty (x^T M x + (u^0)^T R(u^0)) d\tau$$

$$u^0(x) = -R^{-1} B^T \frac{\partial V_u^0(x)}{\partial x} = -R^{-1} B^T P_u^0 x$$

$$V_d^1(x_0) = \max_d \int_0^\infty (x^T M x + (u^0)^T R u^0 - \gamma^2 \|d\|^2) d\tau$$

$$d^1(x) = \frac{1}{\gamma^2} K^T \frac{\partial V_d^1(x)}{\partial x} = \frac{1}{\gamma^2} K^T P_d^1 x$$

智能等级

Level-k: 与已得到的 Level-k - 1 智能进行博弈获得

$$V_u^k(x_0) = \min_u \int_k^\infty (x^T M x + (u^k)^T R(u^k) - \gamma^2 \|d^{k-1}\|^2) d\tau$$

$$u^k(x) = -R^{-1} B^T \frac{\partial V_u^k(x)}{\partial x} = -R^{-1} B^T P_u^k x$$

$$V_d^{k+1}(x_0) = \max_d \int_0^\infty (x^T M x + (u^k)^T R u^k - \gamma^2 \|d^{k+1}\|^2) d\tau$$

$$d^{k+1}(x) = \frac{1}{\gamma^2} K^T \frac{\partial V_d^{k+1}(x)}{\partial x} = \frac{1}{\gamma^2} K^T P_d^{k+1} x$$

Theorem 1

如果满足以下条件，则系统全局渐进稳定

$$\begin{aligned}u^k(0) &= 0, d^{k+1}(0) = 0, \\ \dot{V}_u^k(x) &= -L(x, u^k, d^{k-1}) < 0, \forall x \neq 0, \\ \dot{V}_d^{k+1}(x) &= -L(x, u^k, d^{k+1}) < 0, \forall x \neq 0, \\ H_u^k\left(x, \frac{\partial V_u^k}{\partial x}, u^k, d^{k-1}\right) &= 0, \forall x, \\ H_d^{k+1}\left(x, \frac{\partial V_d^{k+1}}{\partial x} u^k, d^{k+1}\right) &= 0, \forall x, \\ H_u^k\left(x, \frac{\partial V_u^k}{\partial x}, u, d^{k-1}\right) &\geq 0, \forall x, u \\ H_d^{k+1}\left(x, \frac{\partial V_d^{k+1}}{\partial x}, u^k, d\right) &\leq 0, \forall x, d\end{aligned} \tag{5}$$

① 研究背景

② 研究内容

Level- k 模型介绍

零和博弈

认知层分层结构

Model-free Q-learning

③ 计划进度

Q-learning

Q 函数:

$$Q_j^k(x, a_j^k) := V_j^k(x) + H_j^k(x, a_j^k, \nabla V_j^k), \forall x, a_j^k, j \in \{u, d\}$$

增广, 写成矩阵形式:

$$Q_j^k(x, a_j^k) = (U_j^k)^T \begin{bmatrix} Q_{j,xx}^k & Q_{j,xa}^k \\ Q_{j,ax}^k & Q_{j,aa}^k \end{bmatrix} U_j^k := (U_j^k)^T \tilde{Q}_j^k U_j^k$$

Q-learning

当 $j := u$ 时即 $U_u^k = \begin{bmatrix} x^T & (u^k)^T \end{bmatrix}^T$, 此时 \tilde{Q} 矩阵的元素分别为:

$$Q_{u,xx}^k = \left(A + \frac{1}{\gamma^2} K K^T P_d^{k-1} \right)^T P_u^k + P_u^k \left(A + \frac{1}{\gamma^2} K K^T P_d^{k-1} \right) \quad (6)$$

$$+ \left(M - \frac{1}{\gamma^2} P_d^{k-1} K K^T P_d^{k-1} \right) - P_u^k B R^{-1} B^T P_u^k + P_u^k,$$

$$Q_{u,xa}^k = B^T P_u^k, \quad (7)$$

$$Q_{u,ax}^k = B P_u^k, \quad (8)$$

$$Q_{u,aa}^k = R \quad (9)$$

Q-learning

当 $j := d$ 时

即 $U_d^k = \begin{bmatrix} x^T & (d^k)^T \end{bmatrix}^T$, 此时 \tilde{Q} 矩阵的元素分别为:

$$Q_{d,xx}^k = \left(A - BR^{-1}B^T P_u^k \right)^T P_d^k + P_d^k \left(A - BR^{-1}B^T P_u^k \right) \quad (10)$$

$$+ \left(M + P_u^k BR^{-1}B^T P_u^k \right) + \frac{1}{\gamma^2} P_d^k K K^T P_d^k,$$

$$Q_{d,xa}^k = K^T P_d^k, \quad (11)$$

$$Q_{d,ax}^k = K P_d^k, \quad (12)$$

$$Q_{d,aa}^k = -\gamma^2 \quad (13)$$

Q-learning

最优控制

通过庞特里亚金极小原理 $\frac{\partial Q_j^k(x, a_j^k)}{\partial a_j^k} = 0$ 得每个 player 的动作:

$$a_j^k(x) = - \left(Q_{j,aa}^k \right)^{-1} Q_{j,ax}^k, \quad \forall x, j \in \{u, d\}$$

Q-learning

重写 Q 函数

$$Q_j^k(x, a_j^k) = \text{vech}(\tilde{Q})^T (U_j^k \otimes U_j^k), \forall x, a_j^k, j \in \{u, d\}$$

NN 逼近

1、逼近 Q 函数：

$$\hat{Q}_j^k(x, a_j^k) = (\hat{W}_j^k)^T (U_j^k \otimes U_j^k), \quad \forall x, a_j^k, j \in \{u, d\}$$

2、逼近控制 u：

$$\hat{a}_j^k = \hat{W}_{a,j}^k x, \quad \forall x, a_j^k, j \in \{u, d\}$$

积分形式 Q 函数

$$Q_j^k(x(t), a_j^k(t)) = Q_j^k(x(t - T_{\text{IRL}}), a_j^k(t - T_{\text{IRL}})) \quad (14)$$

$$- \int_{t-T_{\text{IRL}}}^t \left(x^T \bar{M}_j^k x + (a_j^k)^T \bar{R}_j a_j^k \right) d\tau, \\ \forall t \geq 0, j \in \{u, d\} \quad (15)$$

其中 $T_{\text{IRL}} \in \mathbb{R}^+$ 是采样时长。

对 Player-1: $\bar{M}_u^k := M - \frac{1}{\gamma^2} P_d^{k-1} K K^T P_d^{k-1}$, $\bar{R}_u := R$;

对 Player-2: $\bar{M}_d^k := M + P_u^k B R^{-1} B^T P_u^k$, $\bar{R}_d := -\gamma^2$

Frame Title

定义：“当前 Q 函数误差项”

$$\begin{aligned}
e_j^k &= \hat{Q}_j^k(x(t), a_j^k(t)) - \hat{Q}_j^k(x(t - T_{\text{IRL}}), a_j^k(t - T_{\text{IRL}})) \\
&\quad + \int_{t-T_{\text{ILL}}}^t \left(x^T \bar{M}_j^k x + (a_j^k)^T \bar{R}_j a_j^k \right) d\tau \\
&= \left(\hat{W}_j^k \right)^T \left(U_j^k(t) \otimes U_j^k(t) \right) - \left(\hat{W}_j^k \right)^T \left(U_j^k(t - T_{\text{IRL}}) \otimes U_j^k(t - T_{\text{ILL}}) \right) \\
&\quad + \int_{t-T_{\text{IRL}}}^t \left(x^T \bar{M}_j^k x + (a_j^k)^T \bar{R}_j a_j^k \right) d\tau,
\end{aligned}$$

定义：“控制策略 u 的误差项”

$$e_{j,a}^k = (W_{a,j}^k)^T x + \left(\hat{Q}_{j,aa}^k \right)^{-1} \hat{Q}_{j,ax} x, \forall x, j \in \{u, d\} \quad (16)$$

Q-learning

更新律

定义 Q 函数的误差: $K_1 = \frac{1}{2} \|e_j^k\|^2$, 和控制策略 u 的误差:

$K_2 = \frac{1}{2} \|e_{j,a}^k\|^2$, 通过梯度下降得到更新律:

$$\dot{\hat{W}}_j^k = -\alpha \frac{\sigma_j^k}{\left(1 + (\sigma_j^k)^T \sigma_j^k\right)^2} \left(e_j^k\right)^T, \forall t \geq 0, j \in \{u, d\} \quad (17)$$

$$\dot{\hat{W}}_{j,a}^k = -\alpha_a \times \left(e_{j,a}^k\right)^T, \forall t \geq 0, j \in \{u, d\} \quad (18)$$

其中 $\sigma_j^k = \left(U_j^k(t) \otimes U_j^k(t)\right) - \left(U_j^k(t - T_{\text{IRL}}) \otimes U_j^k(t - T_{\text{IRL}})\right)$

Lemma 1

误差全局渐进稳定

当增益 α 远远大于增益 α_a 时, 即满足:

$$1 < \alpha_a < \frac{1}{\delta \bar{\lambda} (\bar{R}^{-1})} \left(2\Delta \left(\bar{M}_j^k + Q_{j,xa}^k \bar{R}^{-1} (Q_{j,xa}^k)^T \right) - \bar{\lambda} \left(Q_{j,xa}^k (Q_{j,xa}^k)^T \right) \right)$$

其中 $\Delta = \frac{\sigma_j^k}{(1 + (\sigma_j^k)^T \sigma_j^k)^2}$ 在时间 $[t, t + T_{\text{exp}}]$ 内满足激励条件:

$$\int_t^{t+T_{\text{exp}}} \Delta \Delta^T d\tau \geq \beta I$$

那么 $\psi = \left[x^T \left(\hat{W}_j^k - W_j^k \right)^T \left(\hat{W}_{j,a}^k - W_{j,a}^k \right)^T \right]^T$ 全局渐进稳定

具体证明来自²

²vamvoudakis_q-learning_2017.

① 研究背景

② 研究内容

③ 计划进度

- 研究对象：“Human-Machine” 非零和协作博弈
- 研究目标：safe+stabilisation

Contribution

- 1、Barrier Transformation：建立障碍转换系统，首先保障安全
- 2、Bounded Level-k Rationality：先学习机器和人的 k 个智能等级的行为，以及最优控制行为 ($k=\infty$)
- 3、Human Impact Modeling：将人 (player-1) 的行为利用概率模型进行建模，得到综合行为 (模拟人动作)
- 4、Transfer Learning：基于机器的先验知识，与模拟人的综合行为交互，迭代机器的控制策略 (Online Learning)，实现 stabilisation

Thanks!