

# Project Report

Cooperative Q-Learning for Rejection of  
Persistent Adversarial Inputs in Networked  
Linear Quadratic Systems

By Murtaza Rangwala

ECE6504: Autonomous Coordination

## Overview

This paper addresses the optimization of the distributed performance criteria of agents in a network who are synchronizing to a leader node using cooperative Q-learning. Graphical games based on Game-Theory are generally used to model these multi-agent learning policies. However, most of these game based techniques, require complete modelling of the system, which is not feasible for very large networks. Moreover, solving ODEs in conventional optimal control policies for large networks of multi-agent systems, also suffers from a high computation cost. Q-learning, a popular reinforcement learning algorithm, is applied to this problem, in this paper. Q-learning, is a model-free learning algorithm, which has been applied successfully in this area. However, this learning policy, is well defined for discrete-time based systems, and has not been well established for continuous-time based systems.

The paper proposes a novel way of modelling a continuous time system, and moreover, guarantees convergence to a Nash Equilibrium for a network of agents, synchronizing to one leader. The paper also provides guarantees for an optimal solution, and closed loop stability of the system, while assuming constraints that there is a unique non-negative solution for the saddle point of the Nash Equilibrium.

## Problem Formulation

A strongly connected networked-system,  $G$ , of  $N$  agents is considered, and modelled as follows:

$$\begin{aligned}\dot{x}_i(t) &= A_i x_i(t) + B_i u_i(t) \\ &\quad + D_i v_i(t), \text{ and } x_i(0) \\ &= x_{i0}, t \geq 0, \forall i \in N\end{aligned}$$

Each state,  $x_i(t)$ , is a measurable state vector available for feedback by each agent and the initial conditions are known for all agents.  $u_i(t)$  is the control input, or the minimizing player, and  $v_i(t)$  is adversarial input, or the maximizing player. The assumption of  $(A_i, B_i)$  is stabilizable is made throughout the paper. This is required for the existence of a linear feedback control, Theorem 9.7, Ref [1]. The dynamics of the leader is defined as,  $\dot{x}_L(t) = A_L x_L(t)$ . The agents should seek to cooperatively track the state of the leader node, with the above defined dynamics of the leader. The distributed performance of the system, is then defined as the neighbourhood tracking error of each agent, considering the error of the agents neighbour and its own with the leader.

$$\begin{aligned}e_i &= x_i(t) - x_L(t) \\ e_i &= \sum_{j \in N_i} a_{ij} (x_i - x_j) \\ &\quad + g_i (x_i - x_L), \forall i \in N\end{aligned}$$

For the system to be asymptotically tracking the leader, the neighbourhood error should be,  $e_i \rightarrow 0$ .

The dynamics of the error is then given by,

$$\begin{aligned}\dot{e}_i &= A e_i + (d_i + g_i)(B_i u_i + D_i v_i) - \\ &\quad \sum_{j \in N_i} a_{ij} (B_j u_j - D_j v_j), \forall i \in N\end{aligned}$$

$$\begin{aligned}d_i &= \sum_{j \in N_i} a_{ij}, \text{ defined as the weighted degree of the node } i\end{aligned}$$

The cost function for neighbourhood tracking error, is then defined for each agent. Moreover, the factor of 0.5, is chosen arbitrarily, for a cleaner solution for the following theorems. The weights can be user-defined for each system that you have, and do not affect the optimality and convergence guarantees.

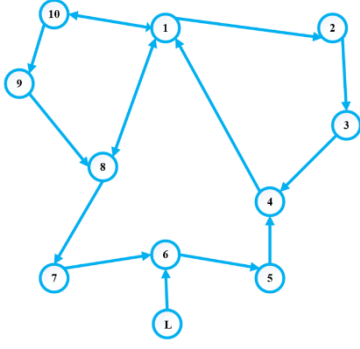


Fig. 1. A networked system  $\mathcal{G}$  with 10 agents and a leader node  $L$ , pinned to agent 6.

The cost function is then defined as,

$$\begin{aligned}
 J_i(e_i(0), u_i, u_{N_i}, v_i, v_{N_i}) &= \frac{1}{2} \int_0^\infty \left( e_i^T H_i e_i \right. \\
 &+ (u_i^T R_{ii} u_i - \gamma_{ii}^2 v_i^T v_i) \\
 &+ \sum_{j \in N_i} (u_j^T R_{ij} u_j \\
 &\left. - \gamma_{ij}^2 v_j^T v_j) \right) dt
 \end{aligned}$$

The problem is generally formulated as an infinite-time horizon linear quadratic soft-constrained differential game. The assumption that,  $(A_i, B_i)$  are stabilizable, and  $(\sqrt{H_i}, A)$  are detectable also guarantees, to ensure a unique non-negative solution, Theorem 4.8 and Theorem 9.7, Ref 1. The detectability of  $(\sqrt{H_i}, A)$ , guarantees that stabilizing control is generated naturally through the solution process, Theorem 9.7.

For a strongly connected graph, with our initial choice of the control input, and adversarial input, we have the graphical Nash equilibrium,

$$\begin{aligned}
 J_i(e_i(0), u_i^*, u_{N_i}^*, v_i^*, v_{N_i}^*) &= \\
 \min_{u_i} J_i(e_i(0), u_i, u_{N_i}^*, v_i^*, v_{N_i}^*) &= \\
 \max_{v_i} J_i(e_i(0), u_i^*, u_{N_i}^*, v_i, v_{N_i}^*), &\quad \text{from} \\
 \text{equation 2.11a.}
 \end{aligned}$$

Using the definition of the value function, from Chapter 4, Ref 1, we get,

$$\begin{aligned}
 V_i^*(e_i(t)) &= \min_{u_i} \max_{v_i} \frac{1}{2} \int_t^\infty \left( e_i^T H_i e_i \right. \\
 &+ (u_i^T R_{ii} u_i - \gamma_{ii}^2 v_i^T v_i) \\
 &+ \sum_{j \in N_i} (u_j^T R_{ij} u_j \\
 &\left. - \gamma_{ij}^2 v_j^T v_j) \right) dt
 \end{aligned}$$

where the weighted norms are considered for the feedback, from time  $t$ .

Note, that the system matrices  $A_i, B_i, D_i$  are not considered for the value function.

## Existence of Nash Equilibrium

Through Theorem 1, proposed in the paper, the authors prove the existence of a Nash Equilibrium given a set of constraints that are satisfied.

The Hamilton-Jacobi-Bellman equation, for the value function that is defined above, using equation 4.60, and Theorems 4.28 and 4.29, is given as:

$$\begin{aligned}
 H_i \left( e_i, u_i, u_{N_i}, v_i, v_{N_i}, \frac{\partial V_i^*}{\partial e_i} \right) &= \\
 &= \frac{\partial V_i^{*T}}{\partial e_i} \left( A e_i + (d_i + g_i)(B_i u_i + D_i v_i) - \sum_{j \in N_i} a_{ij} (B_j u_j - D_j v_j) \right) \\
 &+ \frac{1}{2} \left( e_i^T H_i e_i + (u_i^T R_{ii} u_i - \gamma_{ii}^2 v_i^T v_i) + \sum_{j \in N_i} (u_j^T R_{ij} u_j - \gamma_{ij}^2 v_j^T v_j) \right)
 \end{aligned}$$

Using Theorems 3.2, 3.12, and 3.13, we get the optimal control and worst case adversarial input as

$$u_i^*(e_i) = \arg \min_{u_i} \mathcal{H}_i(e_i, u_i, u_{N_i}, v_i, v_{N_i}, \frac{\partial V_i^*}{\partial e_i})$$

$$= -(d_i + g_i) R_{ii}^{-1} B_i^T \frac{\partial V_i^*}{\partial e_i}, \quad \forall e_i,$$

$$v_i^*(e_i) = \arg \max_{v_i} \mathcal{H}_i(e_i, u_i, u_{N_i}, v_i, v_{N_i}, \frac{\partial V_i^*}{\partial e_i})$$

$$= \frac{(d_i + g_i)}{\gamma_{ii}^2} D_i^T \frac{\partial V_i^*}{\partial e_i}, \quad \forall e_i.$$

The saddle point solution, will satisfy the Hamilton-Jacobi-Bellman equation –

$$H_i \left( e_i, u_i^*, u_{N_i}^*, v_i^*, v_{N_i}^*, \frac{\partial V_i^*}{\partial e_i} \right) = 0$$

Then we have a unique non-negative solution,  $P_i$ , that solves the above equation, where the value function can be represented as,

$$V_i^*(e_i) = \frac{1}{2} e_i^T P_i e_i$$

The above Hamiltonian can then be rewritten as,

$$H_i(e_i, u_i^*, u_{N_i}^*, v_i^*, v_{N_i}^*, e_i^T P_i) = 0$$

$$u_i^*(e_i) = -(d_i + g_i) R_{ii}^{-1} B_i^T P_i e_i, \quad \forall e_i,$$

$$v_i^*(e_i) = \frac{(d_i + g_i)}{\gamma_{ii}^2} D_i^T P_i e_i, \quad \forall e_i.$$

Solving for the above equation, requires a complete knowledge of the system and is also computationally expensive.

**Theorem1:** Let the pairs  $(A, B_i)$  are stabilizable and  $(\sqrt{H_i}, A)$  detectable. Suppose that the graph  $G$  is strongly connected and that there exists unique symmetric positive definite matrices  $P_i$ , that satisfy the coupled Riccati inequalities,

$$H_i(e_i, u_i^*, u_{N_i}^*, v_i^*, v_{N_i}^*, e_i^T P_i) = 0$$

with  $V_i^*(0) = 0, \quad \forall i \in \mathcal{N}$ ,

$$B_i R_{ii}^{-1} B_i^T > \frac{1}{\gamma_{ii}^2} D_i D_i^T, \quad \forall i \in \mathcal{N}$$

$$B_j R_{jj}^{-T} R_{ij} R_{jj}^{-1} B_j^T > \frac{\gamma_{ij}^2}{\gamma_{jj}^4} D_j D_j^T, \quad \forall i, j \in \mathcal{N},$$

with optimal control input,

$$\bar{u}_i(e_i) = -(d_i + g_i) R_{ii}^{-1} B_i^T P_i e_i, \quad \forall e_i,$$

and worst case adversarial input,

$$\bar{v}_i(e_i) = \frac{(d_i + g_i)}{\gamma_{ii}^2} D_i^T P_i e_i, \quad \forall e_i.$$

Then all the agents synchronize asymptotically to the state of the leader, i.e. the equilibrium point of the closed loop system.

$$\dot{V}_i = \frac{\partial V_i}{\partial e_i} \left( A e_i + (d_i + g_i)(B_i \bar{u}_i + D_i \bar{v}_i) - \sum_{j \in \mathcal{N}_i} a_{ij}(B_j \bar{u}_j + D_j \bar{v}_j) \right)$$

$$= -\frac{1}{2} \left( e_i^T H_i e_i + (\bar{u}_i^T R_{ii} \bar{u}_i - \gamma_{ii}^2 \bar{v}_i^T \bar{v}_i) \right.$$

$$\left. + \sum_{j \in \mathcal{N}_i} (\bar{u}_j^T R_{ij} \bar{u}_j - \gamma_{ij}^2 \bar{v}_j^T \bar{v}_j) \right), \quad \forall i \in \mathcal{N},$$

Using eqn. 2.17a as a reference from the textbook, we get an upper bound on the dynamics of the value function,

$$\dot{V}_i \leq -\frac{1}{2} \lambda(H_i) \|e_i\|^2$$

Using  $V_i$  as a Lyapunov function for every agent, the authors prove that,

$$V_i^*(e_i) = 0 \text{ and } \lim_{t \rightarrow \infty} V_i(e_i(t)) = 0$$

With the above conditions, we can arbitrarily define control inputs and adversarial inputs, that satisfy,

$$\mathcal{J}_i(e_i(0); \bar{u}_i, \bar{u}_{N_i}, \bar{v}_i, \bar{v}_{N_i}) = \int_0^\infty \frac{1}{2} \left( e_i^T H_i e_i + (\bar{u}_i^T R_{ii} \bar{u}_i - \gamma_{ii}^2 \bar{v}_i^T \bar{v}_i) \right.$$

$$\left. + \sum_{j \in \mathcal{N}_i} (\bar{u}_j^T R_{ij} \bar{u}_j - \gamma_{ij}^2 \bar{v}_j^T \bar{v}_j) \right) dt$$

$$+ \frac{1}{2} e_i(0)^T P_i e_i(0) + \frac{1}{2} \int_0^\infty \dot{V}_i^* dt$$

$$= \int_0^\infty \frac{1}{2} \left( e_i^T H_i e_i + (\bar{u}_i^T R_{ii} \bar{u}_i - \gamma_{ii}^2 \bar{v}_i^T \bar{v}_i) \right.$$

$$\left. + \sum_{j \in \mathcal{N}_i} (\bar{u}_j^T R_{ij} \bar{u}_j - \gamma_{ij}^2 \bar{v}_j^T \bar{v}_j) \right) dt$$

$$+ \frac{1}{2} e_i(0)^T P_i e_i(0) + \frac{1}{2} \int_0^\infty \frac{\partial V_i^*}{\partial e_i} \left( A e_i + (d_i + g_i)(B_i \bar{u}_i + D_i \bar{v}_i) \right.$$

$$\left. - \sum_{j \in \mathcal{N}_i} a_{ij}(B_j \bar{u}_j + D_j \bar{v}_j) \right) dt, \quad \forall i \in \mathcal{N},$$

$$\begin{aligned} \mathcal{J}_i(e_i(0); \bar{u}_i, u_{\mathcal{N}_i}^*, \bar{v}_i, v_{\mathcal{N}_i}^*) &= \int_0^\infty \frac{1}{2} \left( (\bar{u}_i - u_i^*)^T R_{ii} (\bar{u}_i - u_i^*) \right. \\ &\quad \left. - \gamma_{ii}^2 (\bar{v}_i - v_i^*)^T (\bar{v}_i - v_i^*) \right) dt \\ &\quad + \frac{1}{2} e_i(0)^T P_i e_i(0), \quad \forall i \in \mathcal{N}. \end{aligned}$$

$$\begin{aligned} \mathcal{J}_i(e_i(0); u_i^*, u_{\mathcal{N}_i}^*, \bar{v}_i, v_{\mathcal{N}_i}^*) &= -\frac{1}{2} \gamma_{ii}^2 \int_0^\infty (\bar{v}_i - v_i^*)^T (\bar{v}_i - v_i^*) dt \\ &\quad + \frac{1}{2} e_i(0)^T P_i e_i(0), \quad \forall i \in \mathcal{N}. \end{aligned}$$

$$\begin{aligned} \mathcal{J}_i(e_i(0); \bar{u}_i, u_{\mathcal{N}_i}^*, v_i^*, v_{\mathcal{N}_i}^*) &= \frac{1}{2} \int_0^\infty (\bar{u}_i - u_i^*)^T R_{ii} (\bar{u}_i - u_i^*) dt \\ &\quad + \frac{1}{2} e_i(0)^T P_i e_i(0), \quad \forall i \in \mathcal{N}. \end{aligned}$$

Setting all your inputs as your optimal values, gives

$$\mathcal{J}_i(e_i(0); u_i^*, u_{\mathcal{N}_i}^*, v_i^*, v_{\mathcal{N}_i}^*) = \frac{1}{2} e_i(0)^T P_i e_i(0), \quad \forall i \in \mathcal{N}.$$

This gives us the inequality,

$$\begin{aligned} &Q_i(e_i, u_i, u_{\mathcal{N}_i}, v_i, v_{\mathcal{N}_i}) \\ &= \frac{1}{2} U_i^T \begin{bmatrix} P_i + H_i + P_i A + A^T P_i & (d_i + g_i) P_i B_i & -\text{row}(a_{ij} P_i B_j) & (d_i + g_i) P_i D_i & -\text{row}(a_{ij} P_i D_j)_{j \in \mathcal{N}_i} \\ (d_i + g_i) B_i^T P_i & R_{ii} & 0 & 0 & 0 \\ -\text{col}(a_{ij} B_j^T P_i) & 0 & \text{diag}(R_{ij})_{j \in \mathcal{N}_i} & 0 & 0 \\ (d_i + g_i) D_i^T P_i & 0 & 0 & -\gamma_{ii}^2 & 0 \\ -\text{col}(a_{ij} D_j^T P_i) & 0 & 0 & 0 & -\text{diag}(\gamma_{ij}^2)_{j \in \mathcal{N}_i} \end{bmatrix} U_i \end{aligned}$$

Where,  $U_i =$

$$\begin{bmatrix} e_i^T & u_i^T & u_{\mathcal{N}_i}^T & v_i^T & v_{\mathcal{N}_i}^T \end{bmatrix}^T$$

For the learning rate using the cooperative Q-Learning, we define,

Critic Approximator Weights

$$\begin{aligned} Q_i(e_i, u_i^*, u_{\mathcal{N}_i}^*, v_i^*, v_{\mathcal{N}_i}^*) \\ = W_{ic}^T (U_i \otimes U_i) \end{aligned}$$

$$W_{ic} := \frac{1}{2} \text{vech}(\bar{Q}^i)$$

$$\begin{aligned} &J_i(e_i(0), u_i^*, u_{\mathcal{N}_i}^*, v_i^*, v_{\mathcal{N}_i}^*) \\ &= \min_{u_i} J_i(e_i(0), u_i, u_{\mathcal{N}_i}^*, v_i^*, v_{\mathcal{N}_i}^*) \\ &= \max_{v_i} J_i(e_i(0), u_i^*, u_{\mathcal{N}_i}^*, v_i, v_{\mathcal{N}_i}^*) \end{aligned}$$

## Q-Function

In order to formulate a Q-Function, the authors use the method of defining an Advantage function. This allow the matrix Q, to be state independent, easing computation.

$$\begin{aligned} Q_i(e_i, u_i, u_{\mathcal{N}_i}, v_i, v_{\mathcal{N}_i}) \\ = V_i^*(e_i) \\ + H_i \left( e_i, u_i, u_{\mathcal{N}_i}, v_i, v_{\mathcal{N}_i}, \frac{\partial V_i^*}{\partial e_i} \right) \end{aligned}$$

$$\begin{aligned} &= V_i^*(e_i) + \frac{\partial V_i^*}{\partial e_i}^T \left( A e_i + (d_i + g_i)(B_i u_i + D_i v_i) \right. \\ &\quad \left. - \sum_{j \in \mathcal{N}_i} a_{ij}(B_j u_j + D_j v_j) \right) \\ &\quad + \frac{1}{2} \left( e_i^T H_i e_i + (u_i^T R_{ii} u_i - \gamma_{ii}^2 v_i^T v_i) \right. \\ &\quad \left. + \sum_{j \in \mathcal{N}_i} (u_j^T R_{ij} u_j - \gamma_{ij}^2 v_j^T v_j) \right), \quad \forall e_i, u_i, u_{\mathcal{N}_i}, v_i, v_{\mathcal{N}_i}, \quad \forall i \in \mathcal{N}, \end{aligned}$$

Control Actor Weights

$$\begin{aligned} u_i^*(e_i) &= \min_{u_i} Q_i(e_i, u_i, u_{\mathcal{N}_i}, v_i, v_{\mathcal{N}_i}) \\ &= -(Q_{u_i u_i})^{-1} (Q_{u_i e_i}) e_i \\ &= W_{ia}^T e_i \end{aligned}$$

Adversarial Actor Weights

$$\begin{aligned} v_i^*(e_i) &= \max_{v_i} Q_i(e_i, u_i, u_{\mathcal{N}_i}, v_i, v_{\mathcal{N}_i}) \\ &= -(Q_{v_i v_i})^{-1} (Q_{v_i e_i}) e_i \\ &= W_{id}^T e_i \end{aligned}$$

Consider a system, which requires time 'T', to get excited to read the initial measurements. This system, has its value function defined though a Bellman equation, which writes the value of a decision problem at a certain point in time in terms of the payoff from some initial

choices and the value of the remaining decision problem that results from those initial choices.

$$V_i^*(e_i(t)) = V_i^*(e_i(t-T))$$

$$-\frac{1}{2} \int_{t-T}^t (e_i^T H_i e_i + (u_i^{*T} R_{ii} u_i^* - \gamma_{ii}^2 v_i^{*T} v_i^*)) + \sum_{j \in \mathcal{N}_i} (u_j^{*T} R_{ij} u_j^* - \gamma_{ij}^2 v_j^{*T} v_j^*) d\tau, \forall i \in \mathcal{N},$$

Notice that we can rewrite, the above function as an optimal Q function,

$$E_i = \widehat{W}_{ic}(U_i \otimes U_i) + \frac{1}{2} \int_{t-T}^t \left( e_i^T H_i e_i + (\hat{u}_i^T R_{ii} \hat{u}_i - \gamma_{ii}^2 \hat{v}_i^T \hat{v}_i) + \sum_{j \in \mathcal{N}_i} (\hat{u}_j^T R_{ij} \hat{u}_j - \gamma_{ij}^2 \hat{v}_j^T \hat{v}_j) \right) dt - \widehat{W}_{ic}(U_i(t-T) \otimes U_i(t-T))$$

$$E_{ia} = \widehat{W}_{ia}^T e_i + (\hat{Q}_{u_i u_i}^i)^{-1} (\hat{Q}_{u_i e_i}^i) e_i$$

$$E_{id} = \widehat{W}_{id}^T e_i + (\hat{Q}_{v_i v_i}^i)^{-1} (\hat{Q}_{v_i e_i}^i) e_i$$

We defined the squared-norm errors for each error function, and the gradients are calculated using the least squared method.

$$K_{ic} = \frac{1}{2} \|E_{ic}\|^2$$

$$\dot{\widehat{W}}_{ic} = -\alpha_{ic} \frac{dK_{ic}}{d\widehat{W}_{ic}}$$

Simplifying for the critic estimate gradient, we get,

$$Q_i^*(e_i(t), u_i^*(t), u_{\mathcal{N}_i}^*(t), v_i^*(t), v_{\mathcal{N}_i}^*(t)) = Q^*(e_i(t-T), u_i^*(t-T), u_{\mathcal{N}_i}^*(t-T), v_i^*(t-T), v_{\mathcal{N}_i}^*(t-T)) - \frac{1}{2} \int_{t-T}^t (e_i^T H_i e_i + (u_i^{*T} R_{ii} u_i^* - \gamma_{ii}^2 v_i^{*T} v_i^*)) + \sum_{j \in \mathcal{N}_i} (u_j^{*T} R_{ij} u_j^* - \gamma_{ij}^2 v_j^{*T} v_j^*) d\tau, \forall i \in \mathcal{N}.$$

Now the estimates of the weights for the critic, and the two actors, need to converge to an optimal value. In order to do that, the weights are propagated using gradient descent. An error function of the Q function, and the two actors are calculated for defining the gradient descent. This is nothing but the difference between the LHS and RHS of your above equation.

$$K_{ia} = \frac{1}{2} \|E_{ia}\|^2$$

$$\dot{\widehat{W}}_{ia} = -\alpha_{ia} \frac{dK_{ia}}{d\widehat{W}_{ia}}$$

$$K_{id} = \frac{1}{2} \|E_{id}\|^2$$

$$\dot{\widehat{W}}_{id} = -\alpha_{id} \frac{dK_{id}}{d\widehat{W}_{id}}$$

---

**Algorithm 1: Proposed Q-learning**

---

- 1: **procedure**
  - 2: Start with initial conditions for every agent  $x_i(0)$ , and random initial weights  $\widehat{W}_{ic}(0), \widehat{W}_{ia}(0), \widehat{W}_{id}(0)$  for the critic, actor and worst case adversarial approximator for each agent.
  - 3: Propagate  $t, x_i(t)$  and  $x_j(t)$  in the neighborhood to compute  $e_i$  from (2).
  - 4: Propagate  $\dot{\widehat{W}}_{ic}(t), \dot{\widehat{W}}_{ia}(t), \dot{\widehat{W}}_{id}(t) \triangleright \dot{\widehat{W}}_{ic}$  as in (39),  $\dot{\widehat{W}}_{ia}$  as in (40) and  $\dot{\widehat{W}}_{id}$  as in (41).
  - 5: Compute the Q function  $\hat{Q}_i$  from (31), the optimal control  $\hat{u}_i$  from (32) and the worst case adversarial input  $\hat{v}_i$  from (33) for each agent.
  - 6: **end procedure**
-

$$\begin{aligned}\dot{W}_{ic} &= -\alpha_{ic} \frac{\partial K_{i1}}{\partial W_{ic}} \\ &= -\alpha_{ic} \frac{\left( U_i(t) \otimes U_i(t) - U_i(t-T) \otimes U_i(t-T) \right)}{\left( 1 + \left( U_i(t) \otimes U_i(t) - U_i(t-T) \otimes U_i(t-T) \right)^T \left( U_i(t) \otimes U_i(t) - U_i(t-T) \otimes U_i(t-T) \right) \right)^2} E_i, \quad \forall i \in \mathcal{N},\end{aligned}$$

Note, the gradient is normalized here. The derivation of the gradient, only gives you the numerator. A normalization is added to the gradient.

$$\dot{W}_{ia} = -\alpha_{ia} e_i e_i^T \tilde{W}_{ia} - \alpha_{ia} e_i e_i^T \tilde{Q}_{e_i v_i}^i (\hat{Q}_{u_i u_i}^i)^{-1}, \quad \forall i \in \mathcal{N},$$

$$\dot{W}_{id} = -\alpha_{id} e_i e_i^T \tilde{W}_{id} - \alpha_{id} e_i e_i^T \tilde{Q}_{e_i v_i}^i (\hat{Q}_{v_i v_i}^i)^{-1}, \quad \forall i \in \mathcal{N},$$

Since we know the values for  $\hat{Q}_{v_i v_i}^i$  and  $\hat{Q}_{u_i u_i}^i$  from the look up table, we can simplify the above equation as:

$$\dot{W}_{ia} = -\alpha_{ia} e_i e_i^T \tilde{W}_{ia} - \alpha_{ia} e_i e_i^T \tilde{Q}_{e_i u_i}^i R_{ii}^{-1}, \quad \forall i \in \mathcal{N},$$

$$\dot{W}_{id} = -\alpha_{id} e_i e_i^T \tilde{W}_{id} - \frac{\alpha_{id}}{\gamma_{ii}^2} e_i e_i^T \tilde{Q}_{e_i v_i}^i, \quad \forall i \in \mathcal{N}.$$

**Lemma 2** of the paper proves that the error dynamics of the critic, with fixed control inputs, then the critic weights will have an exponentially stable equilibrium point that satisfies,

$$\begin{aligned}\|\tilde{W}_{ic}(t)\| &\leq \|\tilde{W}_{ic}(t_0)\| \rho_{i1} e^{-\rho_{i2}(t-t_0)}, \quad \forall t > t_0 \\ &\geq 0, \rho_{i1}, \rho_{i2} \in \mathbb{R}^+\end{aligned}$$

This is proved by defining a Lyapunov function

$$\mathcal{L}_i = \frac{1}{2\alpha_{ic}} \|\tilde{W}_{ic}\|^2, \quad t \geq 0, \quad \forall i \in \mathcal{N}.$$

$$\dot{\mathcal{L}}_i = \frac{1}{\alpha_{ic}} \tilde{W}_{ic}^T \dot{\tilde{W}}_{ic}, \quad \forall i \in \mathcal{N}$$

After substituting the value of your dynamics of the estimate function of critic, defined at the top of this page, we can rewrite the equation to,

$$\dot{\mathcal{L}}_i = -\tilde{W}_{ic}^T \Delta_i \Delta_i^T \tilde{W}_{ic}, \quad \forall i \in \mathcal{N}.$$

Where,

$$\Delta_i(t) =$$

$$\frac{\left( U_i(t) \otimes U_i(t) - U_i(t-T) \otimes U_i(t-T) \right)}{1 + \left( U_i(t) \otimes U_i(t) - U_i(t-T) \otimes U_i(t-T) \right)^T \left( U_i(t) \otimes U_i(t) - U_i(t-T) \otimes U_i(t-T) \right)}$$

which is persistently exciting over the complete time interval.

The solution to linear time-varying system, with a state transition matrix defined as  $\frac{\partial \Phi_i(t, t_0)}{\partial t} := -\alpha_{ic} \Delta_i(t) \Delta_i^T(t) \Phi_i(t, t_0)$

$$\tilde{W}_{ic}(t) = \Phi_i(t, t_0) \tilde{W}_{ic}(t_0), \quad \forall t, t_0 > 0, \quad \forall i \in \mathcal{N}$$

The state-transition matrix is exponentially stable, if we have a persistently exciting  $\Delta_i(t)$ .

And for some arbitrary coefficients, we would satisfy,

$$\|\Phi_i(t, t_0)\| = \rho_{i1} e^{-\rho_{i2}(t-t_0)}, \quad \forall t, t_0 > 0, \quad \forall i \in \mathcal{N}.$$

For which we fulfil the criteria,

$$\|\tilde{W}_{ic}(t)\| \leq \|\tilde{W}_{ic}(t_0)\| \rho_{i1} e^{-\rho_{i2}(t-t_0)}, \quad \forall t, t_0 > 0, \quad \forall i \in \mathcal{N}$$

The next step of the paper provides guarantee on the closed loop stability of the origin of  $[e_i^T \quad \tilde{W}_{ic}^T \quad \tilde{W}_{N_i a}^T \quad \tilde{W}_{ia}^T \quad \tilde{W}_{N_i d}^T \quad \tilde{W}_{id}^T]^T$ .

## Theorem 2

The theorem states that considering the neighbourhood tracking error defined as part of this paper, and optimal control and worst case adversarial inputs, the tuning laws generated before and assuming that the tuning law for the critic is faster than the control and adversarial actor, and given the following conditions,

$$\alpha_{ia} > 1; \alpha_{id} > 1,$$

$$\begin{aligned} \alpha_{ia} \bar{\lambda}(R_{ii}^{-1}) + \frac{\alpha_{id}}{\gamma_{ii}^2} &< \frac{2}{\delta_i} \left( \lambda(H_i + Q_{e_i u_i}^i R_{ii}^{-1} (Q_{e_i u_i}^i)^T) \right. \\ &\quad \left. - \frac{1}{2} \bar{\lambda}(Q_{e_i u_i}^i (Q_{e_i u_i}^i)^T + \frac{2}{\gamma_{ii}^2} Q_{e_i v_i}^i (Q_{e_i v_i}^i)^T) \right), \\ \alpha_{ja} \bar{\lambda}(R_{jj}^{-1}) + \frac{\alpha_{jd}}{\gamma_{jj}^2} &< \frac{2}{\delta_j} \left( \lambda(Q_{e_j u_j}^j R_{jj}^{-1} R_{ij} R_{jj}^{-1} (Q_{e_j u_j}^j)^T) \right. \\ &\quad \left. - \frac{1}{2} \bar{\lambda}(Q_{e_j u_j}^j (Q_{e_j u_j}^j)^T + \frac{\gamma_{ij}^2}{\gamma_{jj}^4} Q_{e_j v_j}^j (Q_{e_j v_j}^j)^T) \right), \quad \forall j \in \mathcal{N}_i, \end{aligned}$$

A Lyapunov function, is chosen such that,

$$\begin{aligned} \mathcal{V} = & \sum_{i=1}^N \left( V_i^*(e_i) + \frac{1}{2} \|\tilde{W}_{ic}\|^2 \right. \\ & + \frac{1}{2} \text{trace} \{ \tilde{W}_{ia}^T \tilde{W}_{ia} \} + \frac{1}{2} \sum_{j \in \mathcal{N}_i} \text{trace} \{ \tilde{W}_{ja}^T \tilde{W}_{ja} \} \\ & \left. + \frac{1}{2} \text{trace} \{ \tilde{W}_{id}^T \tilde{W}_{id} \} + \frac{1}{2} \sum_{j \in \mathcal{N}_i} \text{trace} \{ \tilde{W}_{jd}^T \tilde{W}_{jd} \} \right), \end{aligned}$$

Similar, to the dynamics of  $\mathcal{V}$ , derived for above, we can get,

$$\begin{aligned} \dot{\mathcal{V}} = & \sum_{i=1}^N \left( \frac{\partial V_i^*(e_i)}{\partial e_i} \left( A e_i + (d_i + g_i)(B_i \hat{u}_i + D_i \hat{v}_i) \right. \right. \\ & \left. \left. - \sum_{j \in \mathcal{N}} a_{ij}(B_j \hat{u}_j + D_j \hat{v}_j) \right) \right. \\ & + \tilde{W}_{ic}^T \dot{\tilde{W}}_{ic} + \text{trace} \{ \tilde{W}_{ia}^T \dot{\tilde{W}}_{ia} \} + \sum_{j \in \mathcal{N}_i} \text{trace} \{ \tilde{W}_{ja}^T \dot{\tilde{W}}_{ja} \} \\ & \left. + \text{trace} \{ \tilde{W}_{id}^T \dot{\tilde{W}}_{id} \} + \sum_{j \in \mathcal{N}_i} \text{trace} \{ \tilde{W}_{jd}^T \dot{\tilde{W}}_{jd} \} \right), \end{aligned}$$

The above equations can be separated into four different terms,

$$\dot{\mathcal{V}} = \sum_{i=1}^N (T_{i1} + T_{i2} + T_{i3} + T_{i4}),$$

$$\begin{aligned} T_{i1} := & \frac{\partial V_i^*(e_i)}{\partial e_i} \left( A e_i - (d_i + g_i) B_i \tilde{W}_{ia}^T e_i - (d_i + g_i) D_i \tilde{W}_{id}^T e_i \right. \\ & + (d_i + g_i) B_i u_i^* + (d_i + g_i) D_i v_i^* + \sum_{j \in \mathcal{N}} a_{ij} (B_j \tilde{W}_{ja}^T e_j - B_j u_j^*) \\ & \left. + \sum_{j \in \mathcal{N}} a_{ij} (D_j \tilde{W}_{jd}^T e_j - D_j v_j^*) \right), \end{aligned}$$

which is very similar to the Hamilton-Jacobi equations. Since, we know that the Hamiltonian is zero, at the optimal value, we subtract “zero” vis-à-vis, the optimal Hamiltonian from the above equation.

$$\begin{aligned} T_{i1} = & -\frac{1}{2} \sum_{j \in \mathcal{N}_i} (d_j + g_j)^2 e_j^T P_j (B_j R_{jj}^{-T} R_{ij} R_{jj}^{-1} B_j^T \\ & - \frac{\gamma_{ij}^2}{\gamma_{jj}^4} D_j D_j^T) P_j e_j - \frac{1}{2} (d_i + g_i)^2 e_i^T P_i (B_i R_{ii}^{-T} B_i^T \\ & - \frac{1}{\gamma_{ii}^2} D_i D_i^T) P_i e_i - \frac{1}{2} e_i^T H_i e_i \\ & + \frac{\partial V_i^*(e_i)}{\partial e_i} \left( \sum_{j \in \mathcal{N}} a_{ij} (B_j \tilde{W}_{ja}^T + D_j \tilde{W}_{jd}^T) e_j \right. \\ & \left. - (d_i + g_i) (B_i \tilde{W}_{ia}^T + D_i \tilde{W}_{id}^T) e_i \right). \quad ( \end{aligned}$$

$$T_{i3} := -\alpha_{ia} \text{trace} \{ \tilde{W}_{ia}^T e_i e_i^T \tilde{W}_{ia} + \tilde{W}_{ia}^T e_i e_i^T \tilde{Q}_{e_i u_i}^i R_{ii}^{-1} \}$$

$$\begin{aligned} T_{i4} := & -\alpha_{id} \text{trace} \{ \tilde{W}_{id}^T e_i e_i^T \tilde{W}_{id} + \tilde{W}_{id}^T e_i e_i^T \tilde{Q}_{e_i v_i}^i \gamma_{ii}^{-2} \} \\ & - \sum_{j \in \mathcal{N}_i} \alpha_{jd} \text{trace} \{ \tilde{W}_{jd}^T e_j e_j^T \tilde{W}_{jd} + \tilde{W}_{jd}^T e_j e_j^T \tilde{Q}_{e_j v_j}^j \gamma_{jj}^{-2} \}. \end{aligned}$$

The term can be upper bounded after using Young’s inequality,

$$\begin{aligned} T_{i1} \leq & \left( \lambda(H_i + Q_{e_i u_i}^i R_{ii}^{-1} (Q_{e_i u_i}^i)^T) \right. \\ & \left. - \frac{1}{2} \bar{\lambda}(Q_{e_i u_i}^i (Q_{e_i u_i}^i)^T + \frac{2}{\gamma_{ii}^2} Q_{e_i v_i}^i (Q_{e_i v_i}^i)^T) \right) \|e_i\|^2 \\ & - \sum_{j \in \mathcal{N}_i} \left( \lambda(Q_{e_j u_j}^j R_{jj}^{-1} R_{ij} R_{jj}^{-1} (Q_{e_j u_j}^j)^T) \right. \\ & \left. - \frac{1}{2} \bar{\lambda}(Q_{e_j u_j}^j (Q_{e_j u_j}^j)^T + \frac{\gamma_{ij}^2}{\gamma_{jj}^4} Q_{e_j v_j}^j (Q_{e_j v_j}^j)^T) \right) \|e_j\|^2 \\ & + \frac{1}{2} \|e_i^T \tilde{W}_{ia}\|^2 + \frac{1}{2} \sum_{j \in \mathcal{N}_i} \|e_j^T \tilde{W}_{ja}\|^2 \\ & + \frac{1}{2} \|e_i^T \tilde{W}_{id}\|^2 + \frac{1}{2} \sum_{j \in \mathcal{N}_i} \|e_j^T \tilde{W}_{jd}\|^2, \quad ( \end{aligned}$$

$$T_{i2} \leq -\frac{\alpha_{ic}}{4} \|\tilde{W}_{ic}\|^2.$$

$$\begin{aligned} T_{i3} = & -\alpha_{ia} \text{trace} \{ \tilde{W}_{ia}^T e_i e_i^T \tilde{W}_{ia} \} - \alpha_{ia} \text{trace} \{ \tilde{W}_{ia}^T e_i e_i^T \tilde{Q}_{e_i u_i}^i R_{ii}^{-1} \} \\ & - \sum_{j \in \mathcal{N}_i} \alpha_{ja} \text{trace} \{ \tilde{W}_{ja}^T e_j e_j^T \tilde{W}_{ja} + \tilde{W}_{ja}^T e_j e_j^T \tilde{Q}_{e_j u_j}^j R_{jj}^{-1} \} \\ \leq & -\frac{\alpha_{ia}}{2} \|e_i^T \tilde{W}_{ia}\|^2 + \frac{\alpha_{ia} \delta_i}{2} \bar{\lambda}(R_{ii}^{-1}) \|e_i\|^2 \\ & - \sum_{j \in \mathcal{N}_i} \frac{\alpha_{ja}}{2} \|e_j^T \tilde{W}_{ja}\|^2 + \sum_{j \in \mathcal{N}_i} \frac{\alpha_{ja} \delta_j}{2} \bar{\lambda}(R_{jj}^{-1}) \|e_j\|^2, \quad (63) \end{aligned}$$

$$\begin{aligned} T_{i4} = & -\alpha_{id} \text{trace} \{ \tilde{W}_{id}^T e_i e_i^T \tilde{W}_{id} \} - \alpha_{id} \text{trace} \{ \tilde{W}_{id}^T e_i e_i^T \tilde{Q}_{e_i v_i}^i \gamma_{ii}^{-2} \} \\ & - \sum_{j \in \mathcal{N}_i} \alpha_{jd} \text{trace} \{ \tilde{W}_{jd}^T e_j e_j^T \tilde{W}_{jd} + \tilde{W}_{jd}^T e_j e_j^T \tilde{Q}_{e_j v_j}^j \gamma_{jj}^{-2} \} \\ \leq & -\frac{\alpha_{id}}{2} \|e_i^T \tilde{W}_{id}\|^2 + \frac{\alpha_{id} \delta_i}{2 \gamma_{ii}^2} \|e_i\|^2 - \sum_{j \in \mathcal{N}_i} \frac{\alpha_{jd}}{2} \|e_j^T \tilde{W}_{jd}\|^2 \\ & + \sum_{j \in \mathcal{N}_i} \frac{\alpha_{jd} \delta_j}{2 \gamma_{jj}^2} \|e_j\|^2. \end{aligned}$$

$$T_{i2} := -\alpha_{ic} \tilde{W}_{ic}^T \frac{\left( U_i(t) \otimes U_i(t) - U_i(t-T) \otimes U_i(t-T) \right) \left( U_i(t) \otimes U_i(t) - U_i(t-T) \otimes U_i(t-T) \right)^T}{\left( 1 + \left( U_i(t) \otimes U_i(t) - U_i(t-T) \otimes U_i(t-T) \right)^T \left( U_i(t) \otimes U_i(t) - U_i(t-T) \otimes U_i(t-T) \right) \right)^2} \tilde{W}_{ic},$$

Since all of these terms are bounded, when you combine the terms,

$$\begin{aligned} \dot{\mathcal{L}} \leq & \sum_{i=1}^N \left( - \left( \Delta(H_i + Q_{e_i u_i}^i R_{ii}^{-1} (Q_{e_i u_i}^i)^T) \right. \right. \\ & - \frac{1}{2} \bar{\lambda} (Q_{e_i u_i}^i (Q_{e_i u_i}^i)^T + \frac{2}{\gamma_{ii}^2} Q_{e_i v_i}^i (Q_{e_i v_i}^i)^T) \\ & \quad \left. - \frac{\alpha_{ia} \delta_i}{2} \bar{\lambda} (R_{ii}^{-1}) - \frac{\alpha_{id} \delta_i}{2 \gamma_{ii}^2} \right) \|e_i\|^2 \\ & - \sum_{j \in \mathcal{N}_i} \left( \Delta(Q_{e_i u_j}^i R_{jj}^{-1} R_{ij} R_{jj}^{-1} (Q_{e_i u_j}^i)^T) \right. \\ & - \frac{1}{2} \bar{\lambda} (Q_{e_i u_j}^i (Q_{e_i u_j}^i)^T + \frac{\gamma_{ij}^2}{\gamma_{jj}^4} Q_{e_i v_j}^i (Q_{e_i v_j}^i)^T) \\ & \quad \left. - \frac{\alpha_{ja} \delta_j}{2} \bar{\lambda} (R_{jj}^{-1}) - \frac{\alpha_{jd} \delta_j}{2 \gamma_{jj}^2} \right) \|e_j\|^2 \\ & - \frac{\alpha_{ic}}{4} \|\tilde{W}_{ic}\|^2 - \frac{1}{2} (\alpha_{ia} - 1) \|e_i^T \tilde{W}_{ia}\|^2 \\ & - \frac{1}{2} \sum_{j \in \mathcal{N}_i} (\alpha_{ja} - 1) \|e_j^T \tilde{W}_{ja}\|^2 \\ & - \frac{1}{2} (\alpha_{id} - 1) \|e_i^T \tilde{W}_{id}\|^2 - \frac{1}{2} \sum_{j \in \mathcal{N}_i} (\alpha_{jd} - 1) \|e_j^T \tilde{W}_{jd}\|^2 \Big), \end{aligned}$$

We see that the Lyapunov  $s$  is bounded and negative, and hence the closed-loop stability is proved.

## References

1.  $H_\infty$ -Optimal Control and Related Minimax Design Problems, 2<sup>nd</sup> Edition, Basar, Pierre Bernhard

## Commentary on the Derivations

I have understood the complete paper and attempted the derivation of each and every step. The only thing that I haven't understood is the application of the Young's inequality for finding the upper bound on the  $Ti1$  in Theorem 2.