

Formation Control using $GQ(\lambda)$ Reinforcement Learning

by M. Knopp, C. Aykın, J. Feldmaier and H. Shen



Tahsincan Köse
2188423



ORTA DOĞU TEKNİK ÜNİVERSİTESİ
MIDDLE EAST TECHNICAL UNIVERSITY

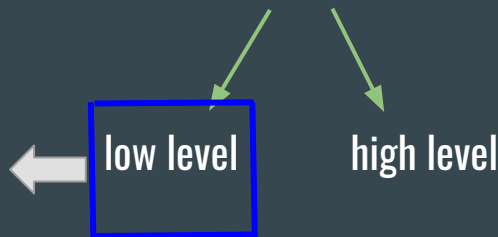
Introduction

Problem: Maintaining a specific shape during locomotion.

- Current solutions in the literature are model, potential-field or graph based approaches which yield **non-linear** control laws.
- Most RL algorithms rely on linear functions which are computationally less complex compared to aforementioned control laws.

Solution: Using RL for Formation Control.

This paper investigates the applicability of the proposed solution on the low level side.



Difficult to implement on low-cost hardware

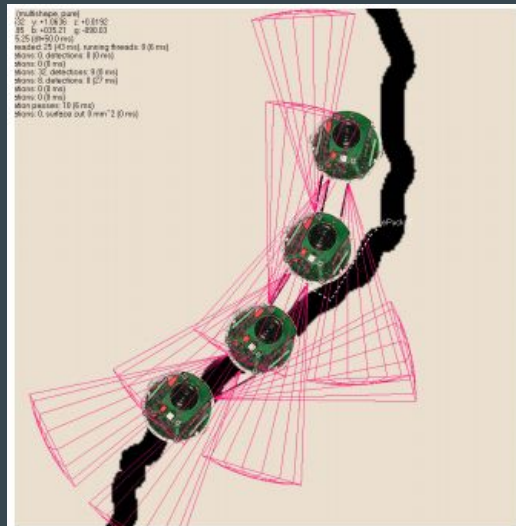


Fig. 3. Simple four agent line formation in V-REP: the first leader can be found in the lower left corner. The red cones are showing active proximity sensors.

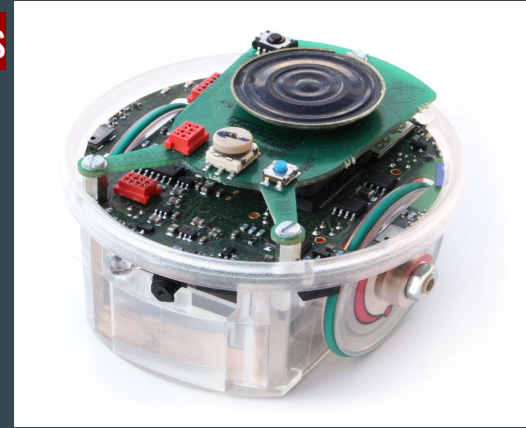
Scenario and Tools

Scenario: Simple leader-follower formation which uses **proximity sensors** as input and motor speeds as actions.

Tools:

- $GQ(\lambda)$ \rightarrow state of art reinforcement learning algorithm.
- ϵ -greedy behavior policy \rightarrow balances exploration and exploitation.
- E-puck robots \rightarrow open-source tabletop robots.
- Gumstix Overo COM \rightarrow Runs embedded Linux and is able to connect to wireless. Used to implement learning algorithms directly on e-puck via Python API.
- V-REP \rightarrow Open-source robot simulation framework.

Maximum range approximately 15 cm



ε -greedy GQ(λ) Learning Algorithm

- $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$; target policy to be learned. Incidentally, if π is chosen as the greedy policy with respect to the learned value function, then the algorithm will implement a generalization of the Greedy-GQ algorithm (Maei, Szepesvari, Bhatnagar & Sutton 2010).
- $\gamma : \mathcal{S} \rightarrow [0, 1]$; termination or discounting function ($\gamma(s) = 1 - \beta(s)$ in GQ paper)
- $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$; reward function

- $b : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$; behavior policy
- $I : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$; interest function (can set to 1 for all state-action pairs or indicate selected state-action pairs to be best approximated)
- $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^n$; feature-vector function
- $\lambda : \mathcal{S} \rightarrow [0, 1]$; bootstrapping or eligibility-trace decay-rate function

Question Functions

Answer Functions

Recall:

$$\pi(s) = \begin{cases} \text{random action from } \mathcal{A}(s) & \text{if } \xi < \varepsilon \\ \operatorname{argmax}_{a \in \mathcal{A}(s)} Q(s, a) & \text{otherwise,} \end{cases}$$

ε -greedy exploration/exploitation



ε -greedy GQ(λ) Learning Algorithm - Continued

- $\theta \in \mathbb{R}^n$; the learned weights of the linear approximation: $Q^\pi(s, a) = \theta^\top \phi(s, a) = \sum_{i=1}^n \theta_i \phi_i(s, a)$
- $w \in \mathbb{R}^n$; secondary set of learned weights
- $e \in \mathbb{R}^n$; eligibility trace vector

Parameters internal to GQ:

- α ; step-size parameter for learning θ
- $\eta \in [0, 1]$; relative step-size parameter for learning w ($\alpha\eta$)

Internal Parameters

$$\rho_t = \frac{\pi(S_t, A_t)}{b(S_t, A_t)}, \quad (1)$$

Importance Sampling Ratio

$$\bar{\phi}_t = \sum_{a \in \mathcal{A}} \pi(S_t, a) \phi(S_t, a) \quad (2)$$

Expected Next Feature Vector



ϵ -greedy GQ(λ) Learning Algorithm - Continued

Goal is to learn this parameter vector.

Full Specification

Pseudocode

$$\delta_t = r(S_t, A_t, S_{t+1}) + \gamma(S_{t+1}) \theta_t^\top \bar{\phi}_{t+1} - \theta_t^\top \phi(S_t, A_t) \quad (3)$$

$$\theta_{t+1} = \theta_t + \alpha \left[\delta_t e_t - \gamma(S_{t+1})(1 - \lambda(S_{t+1}))(\mathbf{w}_t^\top e_t) \bar{\phi}_{t+1} \right] \quad (4)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha \eta [\delta_t e_t - (\mathbf{w}_t^\top \phi(S_t, A_t)) \phi(S_t, A_t)] \quad (5)$$

$$e_t = I(S_t) \phi(S_t, A_t) + \gamma(S_t) \lambda(S_t) \rho_t e_{t-1} \quad (6)$$

weight doubling

The main idea of the algorithm is to minimize the projected Bellman error

$$J(\theta) = \|\Pi T^{\pi_\theta} Q_\theta - Q_\theta\|_\mu^2$$

```

Initialize  $\theta$  arbitrarily and  $\mathbf{w} = 0$ 
Repeat (for each episode):
    Initialize  $e = 0$ 
     $S \leftarrow$  initial state of episode
    Repeat (for each step of episode):
         $A \leftarrow$  action selected by policy  $b$  in state  $S$ 
        Take action  $A$ , observe next state,  $S'$ 
         $\bar{\phi} \leftarrow 0$ 
        For all  $a \in \mathcal{A}(s)$ :
             $\bar{\phi} \leftarrow \bar{\phi} + \pi(S', a) \phi(S', a)$ 
         $\rho = \frac{\pi(S, A)}{b(S, A)}$ 
        GQLearn( $\phi(S, A), \bar{\phi}, \lambda(S'), \gamma(S'), r(S, A, S'), \rho, I(S)$ )
         $S \leftarrow S'$ 
    until  $S'$  is terminal
    
```

```

GQLearn( $\phi, \bar{\phi}, \lambda, \gamma, R, \rho, I$ )
     $\delta \leftarrow R + \gamma \theta^\top \bar{\phi} - \theta^\top \phi$ 
     $e \leftarrow \rho e + I \phi$ 
     $\theta \leftarrow \theta + \alpha(\delta e - \gamma(1 - \lambda)(\mathbf{w}^\top e) \bar{\phi})$ 
     $\mathbf{w} \leftarrow \mathbf{w} + \alpha \eta (\delta e - (\mathbf{w}^\top \phi) \phi)$ 
     $e \leftarrow \gamma \lambda e$ 
    
```



Design of State and Action Space

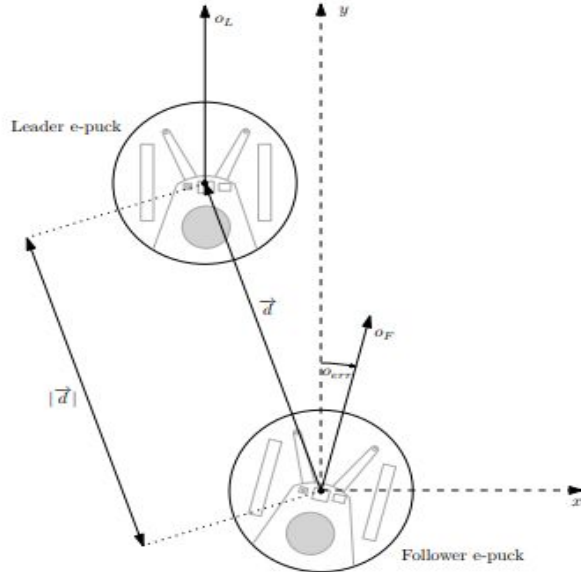


Fig. 1. Illustration of the error measures $|d|$ and o_{err} for our leader-follower scenario.

Since there are two motors in each robot, action space is directly 2-dimensional which represents motor speeds. It is tile coded [linear uniform] as 7×7 for discretization.

State space is determined to be 2-dimensional as well:

- 1) Deviation from optimal distance*, i.e. distance err
- 2) Angular error.

*optimal distance: middle between the maximum and minimum allowed distance

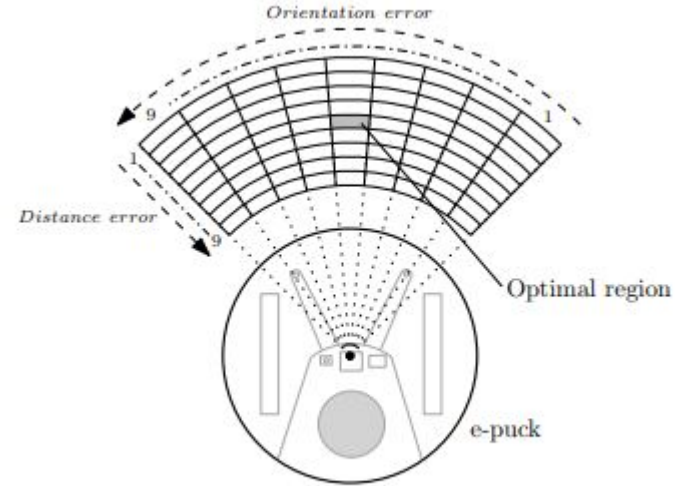


Fig. 2. The 9×9 tile coding used for the representation of the state space when error calculations are based on absolute coordinates.

State space is tile coded [logarithmic] too, as 9×9 for discretization.



Experiments & Results

Simulation

Leader: Simple statistical line follower. Moves arbitrarily.

- Up to 4 followers. Line-like formation. Try to learn the robot directly in front of them.
- Finding: Robots start to stabilize their behavior from first follower towards end of the line formation.

In 74% of all cases, the leader was detected in the dark red area right to optimal region. Additionally, the leader stayed within the orange/red 3×3 area in 98.7% of the time. Logarithmic tiling is used to have rarely visited states distinguishable.

Shift to the right side in the distribution is caused by the overall path which is a closed circular loop.

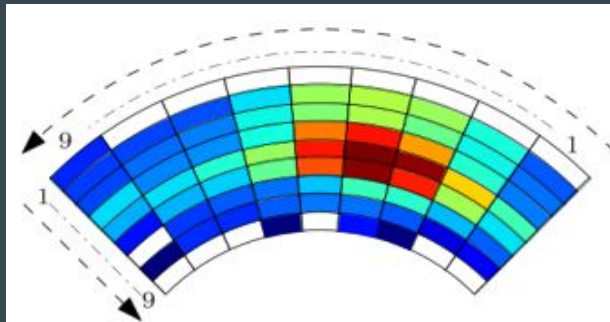


Fig. 4. Distribution of visited states at around 1.3 million steps. In 74% of all cases, the leader was detected in the dark red area right to our optimal region. It stayed within the orange/red 3×3 area on the right side of the optimal region in 98.7% of the time. To make the rarely visited states discernible, a logarithmic scale is used: red and orange values are in the range of 30%–1%, green values are at around 0.1%, and blue values are another magnitude below that. The non-colored states were never visited in this experiment.



Experiments & Results - Continued



Real World

Overo COM instead of a Standart PC is used for whole training.

- Computational bottlenecks are introduced.
- Tiling is reduced from 9x9, 7x7 to 5x5 and 3x3 respectively for state and action values.
- Environmental lighting may heavily distort the sensor readings.
- Tradeoff between smoothness & learning efficiency vs. insurance of followers can keep up with the leader in limited time.

Cure: Transferability of learned parameters from simulation.



Despite all sensor and computational issues, learning algorithm showed that it is able to learn the right parameters for following the leader, i.e. keeping the formation valid.



ORTA DOĞU TEKNİK ÜNİVERSİTESİ
MIDDLE EAST TECHNICAL UNIVERSITY



Nonetheless, stabilization take much more time than the simulation.



Conclusion

- ε -greedy GQ(λ) can be used to have robots learn to keep a line-like formation.
- Required learning parameters can be adapted from a simulation (through transfer of value function estimation).
- “Natural” behavior is obtained through smoothed variant of the leader’s path which would intuitively increase the acceptance of the robot’s behavior by non-experts.

