



Motion planning of a quadrotor robot game using a simulation-based projected policy iteration method*

Li-dong ZHANG^{†1,2}, Ban WANG², Zhi-xiang LIU², You-min ZHANG^{†‡2}, Jian-liang AI¹

¹Department of Aeronautics and Astronautics, Fudan University, Shanghai 200433, China

²Department of Mechanical, Industrial and Aerospace Engineering, Concordia University, Quebec H3G1M8, Canada

[†]E-mail: lidongzhang13@fudan.edu.cn; ymzhang@encs.concordia.ca

Received Sept. 15, 2018; Revision accepted Nov. 27, 2018; Crosschecked Apr. 28, 2019

Abstract: Making rational decisions for sequential decision problems in complex environments has been challenging researchers in various fields for decades. Such problems consist of state transition dynamics, stochastic uncertainties, long-term utilities, and other factors that assemble high barriers including the curse of dimensionality. Recently, the state-of-the-art algorithms in reinforcement learning studies have been developed, providing a strong potential to efficiently break the barriers and make it possible to deal with complex and practical decision problems with decent performance. We propose a formulation of a velocity varying one-on-one quadrotor robot game problem in the three-dimensional space and an approximate dynamic programming approach using a projected policy iteration method for learning the utilities of game states and improving motion policies. In addition, a simulation-based iterative scheme is employed to overcome the curse of dimensionality. Simulation results demonstrate that the proposed decision strategy can generate effective and efficient motion policies that can contend with the opponent quadrotor and gather advantaged status during the game. Flight experiments, which are conducted in the Networked Autonomous Vehicles (NAV) Lab at the Concordia University, have further validated the performance of the proposed decision strategy in the real-time environment.

Key words: Reinforcement learning; Approximate dynamic programming; Decision making; Motion planning; Unmanned aerial vehicle

<https://doi.org/10.1631/FITEE.1800571>

CLC number: TP242

1 Introduction

Autonomous systems, such as unmanned aerial vehicles (UAVs), unmanned ground vehicles (UGVs), and unmanned surface vehicles (USVs) are attracting growing interest in science, industry, and military due to their advantages, such as high maneuverability, extended operational range, and improved personnel

safety, when accomplishing a variety of missions (Yuan et al., 2015; Ghamry et al., 2016; Liu et al., 2016; Sharifi et al., 2016; Wang and Zhang, 2018). Especially along with the rapid development of intelligent systems, autonomous systems have become powerful to deal with difficult and complex tasks. Challenges remain though and are restricted by computational capabilities, sensor reliabilities, communication stabilities, and other limits in both software and hardware. Among these challenges, solving the sequential decision problem plays a critical role in building an autonomous system with high computational efficiency and functional intelligence. A sequential decision problem can be modeled as a Markov decision process (MDP) that

[‡] Corresponding author

* Project supported by the National Natural Science Foundation of China (Nos. 61573282 and 61833013), the Scholarships from China Scholarship Council (No. 201606100139), and the Natural Sciences and Engineering Research Council of Canada

ORCID: You-min ZHANG, <https://orcid.org/0000-0002-9731-5943>

© Zhejiang University and Springer-Verlag GmbH Germany, part of Springer Nature 2019

consists of a set of states, a set of actions at each state, a transition model, and a reward/cost function (Russell and Norvig, 2010). The solution of such a problem, known as a policy, determines which action needs to be taken at each state to achieve the goal state from an initial state. The objective is then to find a feasible and efficient policy that can provide rational decisions for the vehicle in a sophisticated environment to reach the goal criteria. To quantitatively evaluate the performance of the policies, a platform with a certain sequential decision model is needed for further development.

To make rational decisions for quadrotors to reach advantaged status and satisfy goal criteria in the game under a certain formulation, one needs to learn the long-term utilities given sample game states of both quadrotors and terrain information of the game field, and then to generate a near-optimal policy with a long planning horizon. A dynamic programming (DP) algorithm can provide an optimal policy for a sequential decision problem. Nevertheless, it is an intractable problem for an exact DP algorithm to solve a complex game due to the famous curse of dimensionality. McGrew (2008) and McGrew et al. (2010) proposed an approximate dynamic programming (ADP) approach and validated its ability to solve an autonomous air-combat problem. However, they assumed a simplified simulated air-combat problem which restricts the aircraft to level flight with constant velocity. Meanwhile, the state transition dynamics of the air-combat problem were assumed to be a deterministic system. Hence, the size of the game state is exceedingly expanded in comparison to that of state in McGrew et al. (2010). Furthermore, terrain combined with other factors, which was not considered in the previous studies, drains the ability of the aforementioned ADP approach to solve complex decision problems, such as the quadrotor game, using a simple batch manner in the policy evaluation step. In spite of the limited competence of existing schemes, the idea of solving complex sequential decision problems using ADP is still a reasonable one, including establishing the improved efficiency and capability.

Returning to the history of research in sequential decision problems, a fundamental philosophy was established, in which the sequential decision problem could be converted to a problem of determining a one-step optimal operation given the utility

of the remaining decision problem (Bellman, 1952). Therefore, in theory, the optimal solution of a sequential decision problem can be found by a recursive calculation of the Bellman equation until all states are visited. In practice, however, due to the three curses of dimensionality of state space, outcome space, and action space, DP is intractable for real systems (Powell, 2007). Remarkable developments in ADP and reinforcement learning were published and tremendously enhanced the power of DP theory to solve practical problems with large state and action spaces (Bertsekas, 1971; Sutton and Barto, 1998). Their main idea was to approximate utilities of states through linear or nonlinear approximation structures and to generate a suboptimal policy through value iteration (VI) or policy iteration (PI) methods. While VI is more straightforward for application and understanding, the PI method is more efficient, as the approximation of utilities is evaluated policy to policy rather than state to state. Hence, it normally takes few iterations for policies to converge before the convergence of the approximated utility value. There are two general options for approximating the utility value of a certain policy: one is straightforward, and the other is efficient. The first approach, known as the “gradient method,” directly calculates the approximation of the utility value by solving a least squares problem (Bertsekas and Tsitsiklis, 2000; Bertsekas, 2012). The second approach, known as the “indirect method,” calculates the approximated utilities for a fixed policy based on a projected form of the Bellman equation (Yu, 2010, 2012; Bertsekas, 2011). The state-of-the-art research combined the simulation-based methods with the indirect approaches to cut down the dimension of the projected Bellman equation (Buşoniu et al., 2010; Thiery and Scherrer, 2010; Bertsekas, 2015). A multistep version of PI was proposed, where policy improvement was solved based on a multistep Bellman equation (Efroni et al., 2018a,b).

In this study, a one-on-one quadrotor robot game formulation is presented as a platform to examine candidate motion planning strategies. Two deployed quadrotors are assumed to fly inside a cuboid space, containing a static obstacle placed on the ground. The goal for each quadrotor is to maneuver towards the tail of the adversary without being exposed to the front of it. Both quadrotors can move in three dimensions with varying velocities and

independently adjust the facing direction from rotational motions to translational motions because of its unique motion property (Wang et al., 2019). This elegant property enables quadrotors to fly near hovering mode with low speed to make full use of the limited game field space. The proposed formulation assumes that both quadrotors have complete information of the game states. However, input uncertainties are added to the state transition function, of which the state transition probabilities are unknown to the quadrotors. Moreover, the existence of the obstacle introduces additional complexity to the formulation that the utility value of the game states is then determined by not only the relative position and motion status of players but also the relationship of the quadrotors with the terrain. To keep the problem manageable, hard constraints on motions are adopted to avoid collisions with the obstacle. Due to the high flexibility of the quadrotors and the low flight speed, the inflicted constraints are reasonable and adequate for the presented game formulation. To achieve a goal of making rational decisions in the studied quadrotor robot game, a simulation-based projected iterative approach is employed in this study. Specifically, the proposed decision scheme is based on a PI structure, consisting of policy evaluation and policy improvement procedures, where the policy evaluation step is approximately done using a simulation-based projected iterative method. The advantage of this method is that the projected approach is calculated in the feature space rather than the sample space, which greatly increases the ability to handle a complex system with a large size of sample states. The iterative approach effectively avoids large-scale matrix inversions. Accompanying these efforts, the proposed decision scheme can properly work in the quadrotor robot game. The main contributions of this work are the following three aspects: (1) A formulation of a velocity varying one-on-one quadrotor robot game problem in the three-dimensional (3D) space is proposed in conjunction with the stochastic state transition dynamics of the quadrotors; (2) An L_2 regulation method is used to reduce the variance of the policy approximation structure, and hence improve the stability of the generated policy; (3) Real flight tests using the latest Quanser QDrone UAVs are conducted, validating the performance of the proposed motion planning algorithm in the real-time environment.

2 Quadrotor robot game formulation

The motion planning problem of the quadrotor robot game is formulated as a stationary discounted total cost minimization problem of finite states, in which the game state is updated based on the first-order Markov model, expressed as

$$\mathbf{s}_{k+1} = f(\mathbf{s}_k, \mathbf{a}_k, \boldsymbol{\omega}_k), \quad k = 1, 2, \dots, \quad (1)$$

where the state \mathbf{s}_k is an element of a space S , the action \mathbf{a}_k an element of a space A , and the random disturbance $\boldsymbol{\omega}_k$ an element of a space D . S , A , and D are assumed to be finite spaces with countable elements. The player in Eq. (1) always decides to invariably behave at the same state; that is to say, the optimal action at a given state would not change over time. Note that there is no fixed time limit in a stationary problem; thus, it is recognized as an infinite-horizon problem. In this case, if the player does not ever satisfy the goal criteria, the utility with additive and undiscounted costs will generally be infinite, and then it turns out to be difficult to find the optimal policy correlated with the lowest infinite utility value (Russell and Norvig, 2010). Therefore, a discount factor $\gamma \in (0, 1)$ for the cost per stage is introduced, and the cost functions are bounded. The objective is then to find a stationary policy μ , where $\mu(\mathbf{s}_k) \in A(\mathbf{s}_k)$ for all $\mathbf{s}_k \in S$ ($k = 1, 2, \dots$), which minimizes the total cost function as

$$J_\mu(\mathbf{s}_i) = \lim_{N \rightarrow \infty} \mathbb{E} \left\{ \sum_{k=0}^{N-1} \gamma^k g(\mathbf{s}_k, \mathbf{s}_{k+1}) \mid \mathbf{s}_0 = \mathbf{s}_i \right\}, \quad (2)$$

where $g(\mathbf{s}_k, \mathbf{s}_{k+1})$ is the function of cost per stage. By applying DP mapping to J_μ , we can obtain the total cost function in the Bellman equation as

$$\begin{aligned} (T_\mu J_\mu)(\mathbf{s}_i) &= \mathbb{E} \{ g(\mathbf{s}_i, \mathbf{s}_j) + \gamma J_\mu(f(\mathbf{s}_i, \mathbf{a}_i, \boldsymbol{\omega}_i)) \} \\ &= \sum_{j=0}^n p_{ij}(\mu(\mathbf{s}_i))(g(\mathbf{s}_i, \mathbf{s}_j) + \gamma J_\mu(\mathbf{s}_j)), \end{aligned} \quad (3)$$

where p_{ij} ($i, j = 1, 2, \dots, n$) are the state transition probabilities, and n is the size of the state space S .

In this quadrotor robot game formulation, the game state \mathbf{s} is defined by positions, velocities, and facing angles of two players in blue and red, identified as

$$\mathbf{s} = [x_p, y_p, z_p, u_p, v_p, w_p, \psi_p], \quad p = \text{b, r}, \quad (4)$$

where \mathbf{s} is a vector with 14 elements, b represents blue, and r represents red. Constraints are added

to motions to keep the quadrotors within a cuboid game field and to avoid collisions with the obstacle placed on the ground (Fig. 1).

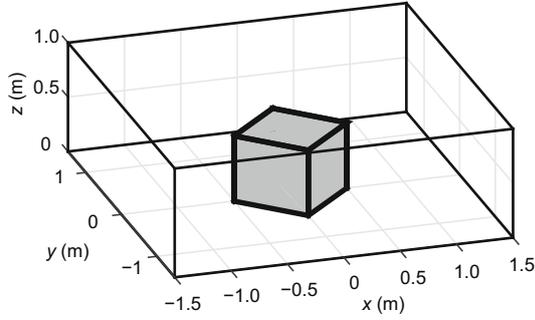


Fig. 1 Game field of the quadrotor robot game

Quadrotors are allowed to fly in three dimensions, and can make pure yaw motion unallied with translational motions. Hence, the action space is set as $A = \{u^+, u^-, v^+, v^-, w^+, w^-, \psi^+, \psi^-, 0\}$, which means that $\mathbf{a}_p = [a_p^u, a_p^v, a_p^w, a_p^\psi] \in \{[1, 0, 0, 0], [-1, 0, 0, 0], [0, 1, 0, 0], [0, -1, 0, 0], [0, 0, 1, 0], [0, 0, -1, 0], [0, 0, 0, 1], [0, 0, 0, -1], [0, 0, 0, 0]\}$, $p = b, r$. At each step ($\Delta t = 0.25$ s), game states are updated by the state transition function (Algorithm 1). Values of parameters used in the state transition function are listed in Table 1. The obstacle interval and threshold height values for a safe pass over the obstacle are selected based on the size of the quadrotor. In the studied game formulation, $\omega^u, \omega^v, \omega^w$, and ω^ψ are selected from the finite set $D = \{-0.2, -0.15, 0, 0.15, 0.2\}$ in accordance with a probability P_ω . However, the actual distribution of P_ω is unknown to the players. At the beginning, an empirical priori probability

Table 1 Parameters of the state transition function

Parameter	Value	Unit	Parameter	Value	Unit
\dot{u}_b	0.2	m/s ²	\dot{u}_r	0.2	m/s ²
\dot{v}_b	0.2	m/s ²	\dot{v}_r	0.2	m/s ²
\dot{w}_b	0.1	m/s ²	\dot{w}_r	0.1	m/s ²
$\dot{\psi}_b$	$\pi/9$	rad/s	$\dot{\psi}_r$	$\pi/9$	rad/s
$u_{b\max}$	0.2	m/s	$u_{r\max}$	0.2	m/s
$v_{b\max}$	0.2	m/s	$v_{r\max}$	0.2	m/s
$w_{b\max}$	0.1	m/s	$w_{r\max}$	0.1	m/s
$u_{b\min}$	-0.2	m/s	$u_{r\min}$	-0.2	m/s
$v_{b\min}$	-0.2	m/s	$v_{r\min}$	-0.2	m/s
$w_{b\min}$	-0.1	m/s	$w_{r\min}$	-0.1	m/s
x_{\max}	1.5	m	x_{\min}	-1.5	m
y_{\max}	1.5	m	y_{\min}	-1.5	m
z_{\max}	1.0	m	z_{\min}	0.2	m

Algorithm 1 State transition function $f(\mathbf{s}, \mathbf{a}, \boldsymbol{\omega})$

```

1: for  $p = b, r$  do
2:   if  $x_p$  in an obstacle interval and  $z_p < z_{\text{safe}}$  then
3:      $u_p = 0$ 
4:   else
5:      $u_p = u_p + \dot{u}_p(a_p^u + \omega^u)\Delta t$ 
6:     // Update the velocity component in the x axis
7:      $u_p = \min(\max(u_p, u_{p\min}), u_{p\max})$ 
8:     // Constraints on u
9:   end if
10:  if  $y_p$  in an obstacle interval and  $z_p < z_{\text{safe}}$  then
11:     $v_p = 0$ 
12:  else
13:     $v_p = v_p + \dot{v}_p(a_p^v + \omega^v)\Delta t$ 
14:    // Update the velocity component in the y axis
15:     $v_p = \min(\max(v_p, v_{p\min}), v_{p\max})$ 
16:    // Constraints on v
17:  end if
18:  if  $z_p$  in an obstacle interval and  $z_p < z_{\text{safe}}$  then
19:     $w_p = 0$ 
20:  else
21:     $w_p = w_p + \dot{w}_p(a_p^w + \omega^w)\Delta t$ 
22:    // Update the velocity component in the z axis
23:     $w_p = \min(\max(w_p, w_{p\min}), w_{p\max})$ 
24:    // Constraints on w
25:  end if
26:   $\psi_p = \psi_p + \dot{\psi}_p(a_p^\psi + \omega^\psi)\Delta t$  // Update the yaw angle
27:   $x_p = x_p + u_p\Delta t$ ;  $x_p = \min(\max(x_p, x_{\min}), x_{\max})$ 
28:   $y_p = y_p + v_p\Delta t$ ;  $y_p = \min(\max(y_p, y_{\min}), y_{\max})$ 
29:   $z_p = z_p + w_p\Delta t$ ;  $z_p = \min(\max(z_p, z_{\min}), z_{\max})$ 
30: end for

```

$P_{\omega 0} = \{0.05, 0.10, 0.70, 0.10, 0.05\}$ is adopted. As the game proceeds, P_ω is updated by the frequencies of visited states. The disturbance set D is assumed to be invariant in spite of the changes of P_ω ; thus, the size of the state space S remains unchanged throughout the game.

Cost per stage function $g(\mathbf{s}_i, \mathbf{s}_j)$ consists of three contributions: orientation, projected range, and terrain. Fig. 2 illustrates the projected geometric relationship of two players on the x - y plane. Since the quadrotors fly near hover mode, orientation angles are thus defined on the x - y plane. The cost function

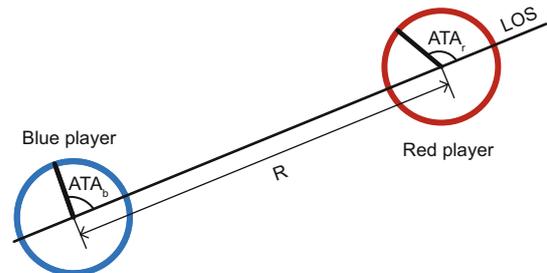


Fig. 2 Projected relative geometry of two players

References to color refer to the online version of this figure

contribution of orientation is calculated as

$$S_A = \frac{|ATA_b| + |ATA_r|}{\pi} - 1, \quad (5)$$

where ATA is the antenna train angle (ATA), measured from the nose of a quadrotor to the line of sight (LOS) vector. LOS refers to the vector connecting the centers of mass from the blue quadrotor to the red quadrotor. When the blue quadrotor is directly facing the red quadrotor, $ATA_b = 0$; if the red quadrotor is directly facing the blue quadrotor, $ATA_r = \pm\pi$. On the other hand, if the blue quadrotor is directly behind the red quadrotor, then $ATA_b = \pm\pi$; vice versa for the red quadrotor. Hence, the blue player is going to choose actions that minimize S_A , while the red player prefers actions that maximize S_A . Since both ATA_b and ATA_r are within a closed interval $[-\pi, \pi]$, S_A is a bounded cost function. Note that because of the ability to make pure yaw motion independent of translational motions, ATA may be different with the angle between the projected velocity on the $x-y$ plane and the LOS vector, which quite commonly happens in a quadrotor game. With the contribution of orientation, quadrotors may take actions to aim at the tail of the opponent; nonetheless, it may try to fly away from the opponent to prevent exposing its back. Therefore, a contribution of range S_R should attract both players to actively engage in the game. When the projected relative range R is larger than the desired range R_d , S_R pulls the quadrotor towards the opponent; if the projected range value R is smaller than a safety threshold range R_{th} , S_R pushes the quadrotor away from the opponent to avoid collisions. In practice, a buffer range is added to R_d to prevent overshoot, thus decreasing the risk of collision with the opponent. The cost function of range is formulated as

$$S_R = \begin{cases} k_{r1}(R - R_d), & R < R_d, \\ 0, & R_d \leq R \leq R_{th}, \\ k_{r2}(R - R_{th}), & R > R_{th}. \end{cases} \quad (6)$$

Note that unlike the air-combat scenario, concepts of pursuing and evading become subtle in the quadrotor game because of the independence of the facing direction and the direction of velocity. Hence, in most circumstances, both players attempt to minimize S_R to achieve the desired range. The cost function combined with orientation and projected range contributions can encourage two players to make maneuvers on the $x-y$ plane to obtain the advantage in

the game with a constant clearance from the ground. However, the various terrains offer possibilities for achieving advantaged status rather than solely flying on the $x-y$ plane. For instance, if the opponent is behind the obstacle, the quadrotor can decide either to make a detour while maintaining the height or to fly across the terrain by ascending above the obstacle. Consequently, the contributions of terrain inside and outside the obstacle area are respectively defined as

$$S_H = \begin{cases} 0, & z_{safe} < z < z_{max}, \\ k_{h0}(z - z_d), & z < z_{safe}, \end{cases} \quad (7)$$

$$S_H = \begin{cases} k_{h1}(z - z_{d+}), & z_{d+} < z \leq z_{max}, \\ 0, & z_{d-} < z < z_{d+}, \\ k_{h2}(z - z_{d-}), & z < z_{d-}, \end{cases} \quad (8)$$

where $z_{d+} = z_d + z_{buffer}$ and $z_{d-} = z_d - z_{buffer}$.

When the quadrotor is outside the obstacle interval or the opponent is flying higher than the obstacle, the desired height z_d is merely the current height value of the opponent; otherwise, z_d equals the safety threshold height to fly over the obstacle. Similar to the contribution of range, both blue and red players attempt to obtain the minimum S_H to achieve the desired height in the game. Moreover, since R and z are within the closed intervals, both S_R and S_H are bounded by the cost functions. Therefore, the costs per stage for blue and red players are defined as

$$\begin{cases} g_b = S_A + S_R + S_{H_b}, \\ g_r = S_A - S_R - S_{H_r}. \end{cases} \quad (9)$$

The cost per stage function can stimulate the quadrotor to win the game in a greedy manner, where the blue player prefers actions that can minimize g_b , and the red player chooses actions that can maximize g_r . However, the decision strategy for a sequential decision making problem can perform better if the decision is made based on a longer horizon, rather than a single step look ahead greedy method. The presented formulation of the quadrotor robot game provides a platform to evaluate the performance of candidate decision algorithms given the same game environment, including game field setup, quadrotor dynamics, and cost per stage functions. In the subsequent section, the proposed motion planning algorithm based on the ADP approach is introduced. This algorithm will then be adopted by the blue player in Section 4 to evaluate its performance by playing a quadrotor robot game against

the red player driven by the *-minimax algorithm (Ballard, 1983; Hauk et al., 2004).

3 Simulation-based projected policy iteration method

Starting from an initial policy, the PI method generates an improved policy by iteratively executing policy evaluation and policy improvement calculations (Fig. 3). In an ADP-based scheme, the policy evaluation step is approximately done by one feasible approximation architecture. At each iteration,

$$\tilde{J}_\mu(\mathbf{s}_i, \mathbf{r}) = \phi^\top(\mathbf{s}_i)\mathbf{r}, \quad (10)$$

$$\mu'(\mathbf{s}_i) = \arg \min_{a \in A} \sum_{j=1}^n p_{ij}(a)(g(\mathbf{s}_i, \mathbf{s}_j) + \gamma \tilde{J}_\mu(\mathbf{s}_j, \mathbf{r})), \quad (11)$$

where the utility of the current policy μ at the state \mathbf{s}_i is approximated by a linear structure consisting of a feature vector $\phi(\mathbf{s}_i) \in \mathbb{R}^q$ and a fit coefficient vector $\mathbf{r} \in \mathbb{R}^q$. An improved policy $\mu'(\mathbf{s}_i)$ is adopted by the one-step Bellman equation. An optimal approximation of \mathbf{J}_μ can be found by solving a least squares optimization problem, defined as

$$\mathbf{r}^* = \arg \min \|\mathbf{J}_\mu - \tilde{\mathbf{J}}_\mu\|_2^2, \quad (12)$$

where $\tilde{\mathbf{J}}_\mu = \Phi \mathbf{r}$ and $\Phi = [\phi^\top(\mathbf{s}_1), \phi^\top(\mathbf{s}_2), \dots, \phi^\top(\mathbf{s}_n)] \in \mathbb{R}^{n \times q}$. Assume that Φ has a rank q . Then \mathbf{r}^* can be directly solved by the inversion in batch manner as

$$\mathbf{r}^* = (\Phi^\top \Phi)^{-1} \Phi^\top \hat{\mathbf{J}}_\mu, \quad (13)$$

where $\hat{\mathbf{J}}_\mu$ is the estimated value of \mathbf{J}_μ based on a sufficient number of sample states. Such an approach has been validated (McGrew, 2008; McGrew et al., 2010; Ma et al., 2014; Fang et al., 2016). However, in the cases where the state space is very large and the number of necessary sample states is also large, calculating Eq. (13) becomes impracticable or extremely time consuming if numerical methods are used. This challenge can be overcome by approximating the utility value \mathbf{J}_μ in the feature space rather than in the state space.

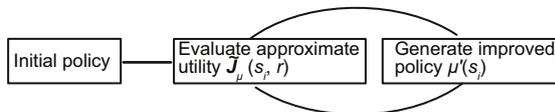


Fig. 3 Policy iteration scheme

To solve this problem, we introduce the projected form of the Bellman equation as

$$\Phi \mathbf{r} = \Pi \mathbf{T}_\mu(\Phi \mathbf{r}), \quad (14)$$

where Π is the projection operator, and \mathbf{T}_μ is the Bellman operator with respect to policy μ . Eq. (14) is an approximation of the Bellman equation $\mathbf{J}_\mu = \mathbf{T}_\mu \mathbf{J}_\mu$ on the feature space. \mathbf{T}_μ is a simple matrix form in Eq. (3), which can be written as

$$\mathbf{T}_\mu \mathbf{J}_\mu = \mathbf{g} + \gamma \mathbf{P} \mathbf{J}_\mu, \quad (15)$$

where $\mathbf{g} = \sum_{j=0}^n p_{ij} g(\mathbf{s}_i, \mathbf{s}_j)$, and \mathbf{P} is the matrix with components p_{ij} . The project operator Π denotes a projection onto the feature space, which minimizes a weighted Euclidean norm $\|\cdot\|_\xi^2$ with positive weights ξ_i ($i = 1, 2, \dots, n$), written as

$$\begin{cases} \Pi \mathbf{J}_\mu = \Phi \mathbf{r}_\Pi, \\ \mathbf{r}_\Pi = \arg \min_{\mathbf{r} \in \mathbb{R}^q} \|\mathbf{J}_\mu - \Phi \mathbf{r}\|_\xi^2 \\ = \arg \min_{\mathbf{r} \in \mathbb{R}^q} \sum_{i=1}^n \xi_i (J_\mu(\mathbf{s}_i) - \phi^\top(\mathbf{s}_i)\mathbf{r})^2. \end{cases} \quad (16)$$

Assume that Φ has a rank q . Then \mathbf{r}_Π can be solved by taking the gradient of $\|\mathbf{J}_\mu - \Phi \mathbf{r}\|_\xi^2$ to 0, expressed as

$$\mathbf{r}_\Pi = (\Phi^\top \Xi \Phi)^{-1} \Phi^\top \Xi \mathbf{J}_\mu. \quad (17)$$

Thus, we have

$$\Pi = \Phi \mathbf{r}_\Pi = \Phi (\Phi^\top \Xi \Phi)^{-1} \Phi^\top \Xi, \quad (18)$$

where Ξ is a diagonal matrix with $\Xi_{ii} = \xi_i$, $i = 1, 2, \dots, n$. Combining the operators Π and \mathbf{T}_μ leads to

$$\Pi \mathbf{T}_\mu(\Phi \mathbf{r}) = \Phi (\Phi^\top \Xi \Phi)^{-1} \Phi^\top \Xi (\mathbf{g} + \gamma \mathbf{P} \Phi \mathbf{r}). \quad (19)$$

It has been proved that if ξ refers to a steady-state probability vector, then $\Pi \mathbf{T}_\mu$ is a contraction with respect to $\|\cdot\|_\xi$. Therefore, the project Bellman equation has a unique solution and can be successively solved by applying $\Pi \mathbf{T}_\mu$, known as the projected value iteration (PVI) method (Bertsekas, 2007). Starting with an arbitrary initial vector $\Phi \mathbf{r}_0$, the iteration is formed as

$$\Phi \mathbf{r}_{k+1} = \Pi \mathbf{T}_\mu(\Phi \mathbf{r}_k). \quad (20)$$

By substituting Eq. (19) into Eq. (20), the iteration can be explicitly written as

$$\mathbf{r}_{k+1} = \mathbf{r}_k - (\Phi^\top \Xi \Phi)^{-1} (\mathbf{C} \mathbf{r}_k - \mathbf{d}), \quad (21)$$

where

$$\begin{cases} \mathbf{C} = \Phi^T \Xi (\mathbf{I} - \gamma \mathbf{P}) \Phi, \\ \mathbf{d} = \Phi^T \Xi \mathbf{g}. \end{cases} \quad (22)$$

\mathbf{r}_k converges to a steady vector \mathbf{r}^* after sufficient iterations, while $\Phi \mathbf{r}^*$ is the solution of the projected Bellman equation, i.e., the optimal approximation of the utility value. Nevertheless, PVI still depends on matrix products of size n , which again becomes impracticable for a large and complex problem. Furthermore, the probability vector ξ is generally unknown at the beginning. These difficulties can be overcome if we treat the weighted Euclidean norm as an expected value with respect to a probability distribution (ξ here is a steady-state probability vector). Therefore, the projection operation can be viewed as a Monte Carlo simulation, which then allows us to use simulation-based methods to evaluate the approximate utility value.

The idea is to express the least squares problem as an expected value with respect to a probability distribution (ξ), and then approximate the expected value by sampling depending on that distribution. For instance, by recalling Eqs. (16) and (17), \mathbf{r}_Π can be written in an element form as

$$\mathbf{r}_\Pi = \left(\sum_{i=1}^n \xi_i \phi(\mathbf{s}_i) \phi^T(\mathbf{s}_i) \right)^{-1} \sum_{i=1}^n \xi_i \phi(\mathbf{s}_i) J_\mu(\mathbf{s}_i). \quad (23)$$

Suppose a sequence of samples of $t = 0, 1, \dots$, are generated regarding the fixed stationary policy μ according distribution ξ . At $t = k$, \mathbf{r}_Π can be estimated as

$$\begin{aligned} \hat{\mathbf{r}}_{\Pi k} &= \left(\frac{1}{k+1} \sum_{t=0}^k \phi(\mathbf{s}_t) \phi^T(\mathbf{s}_t) \right)^{-1} \left(\frac{1}{k+1} \sum_{t=0}^k \phi(\mathbf{s}_t) J_\mu(\mathbf{s}_t) \right) \\ &= \left(\sum_{t=0}^k \phi(\mathbf{s}_t) \phi^T(\mathbf{s}_t) \right)^{-1} \left(\sum_{t=0}^k \phi(\mathbf{s}_t) J_\mu(\mathbf{s}_t) \right), \end{aligned} \quad (24)$$

where k denotes the number of simulation samples. As the simulation is executed in a discrete manner, the simulation pace can be unified with that of the iteration of \mathbf{r} . Hence, \mathbf{r} is updated at each step when a new simulation sample is generated. Note that $\hat{\mathbf{r}}_\Pi$ is guaranteed to converge to \mathbf{r}_Π as the number of simulation samples k increases. To clarify this, we are coping with a finite-state MDP with a fixed policy; thus, the distribution of state \mathbf{s}_i is consistent with the long-term frequencies of occurrence of \mathbf{s}_i

during the simulation,

$$\xi_i = \lim_{k \rightarrow \infty} \frac{1}{k+1} \sum_{t=0}^k \delta(\mathbf{s}_t = \mathbf{s}_i), \quad i = 1, 2, \dots, n. \quad (25)$$

If $\mathbf{s}_t = \mathbf{s}_i$, $\delta(\mathbf{s}_t = \mathbf{s}_i) = 1$; otherwise, $\delta(\mathbf{s}_t = \mathbf{s}_i) = 0$.

By applying the simulation method to PVI, we can reform Eqs. (21) and (22) as

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \mathbf{G}_k^{-1} (\mathbf{C}_k \mathbf{r}_k - \mathbf{d}_k), \quad (26)$$

where

$$\begin{cases} \mathbf{G}_k = \frac{1}{k+1} \sum_{t=0}^k \phi(\mathbf{s}_t) \phi^T(\mathbf{s}_t), \\ \mathbf{C}_k = \frac{1}{k+1} \sum_{t=0}^k \phi(\mathbf{s}_t) (\phi(\mathbf{s}_t) - \gamma \phi(\mathbf{s}_{t+1}))^T, \\ \mathbf{d}_k = \frac{1}{k+1} \sum_{t=0}^k \phi(\mathbf{s}_t) g(\mathbf{s}_t, \mathbf{s}_{t+1}). \end{cases} \quad (27)$$

By substituting Eq. (27) into Eq. (26), the iteration can be written as

$$\begin{aligned} \mathbf{r}_{k+1} &= \left(\sum_{t=0}^k \phi(\mathbf{s}_t) \phi^T(\mathbf{s}_t) \right)^{-1} \\ &\cdot \left(\sum_{t=0}^k \phi(\mathbf{s}_t) (\gamma \phi^T(\mathbf{s}_{t+1}) \mathbf{r}_k + g(\mathbf{s}_t, \mathbf{s}_{t+1})) \right), \end{aligned} \quad (28)$$

where \mathbf{r}_{k+1} eventually converges to the solution of the simulation-based projected Bellman equation \mathbf{r}^* . This simulation-based PVI method is known as the ‘‘least squares policy evaluation’’ (LSPE) method. The key benefit of the LSPE method is that it involves only q -dimensional matrix-vector calculations, and q can be significantly smaller than n . Moreover, the distribution vector ξ is not required to be predetermined. We can now derive the simulation-based projected policy iteration method as

$$\tilde{J}_{\mu,k}(\mathbf{s}_i, \mathbf{r}_k) = \phi^T(\mathbf{s}_i) \mathbf{r}_k, \quad (29)$$

$$\mu'_k(\mathbf{s}_i) = \arg \min_{a \in A} \sum_{j=1}^n p_{ij}(a) (g(\mathbf{s}_i, \mathbf{s}_j) + \gamma \tilde{J}_{\mu,k}(\mathbf{s}_j, \mathbf{r}_k)), \quad (30)$$

where \mathbf{r}_k is updated by Eq. (28). By successively applying Eqs. (29) and (30), a suboptimal policy μ^* can be found for the stationary discounted total cost minimization problem to be solved.

To obtain the solution of the projected Bellman equation, it is assumed that Φ has a rank q ; otherwise, the matrix $\Phi^T \Xi \Phi$ is singular, and thus it is

difficult to calculate the inversion. As the complexity of the problem increases, more features are needed to represent the full properties of game states, and it becomes more challenging to guarantee that all features are linearly independent. Suppose there are two correlated feature vectors in Φ , and this may cause a pair of fit coefficients in \mathbf{r} to have large values with the opposite sign to cancel each other (Hastie et al., 2001). The variance of the approximate structure becomes worse as the correlation enhances, which may eventually cause the instability of the decision process. A less variant structure can be obtained by introducing L_2 regulation to the least squares term in Eq. (16) with a penalty term on \mathbf{r} . Thus, the least squares problem can be reformed as

$$\mathbf{r}_{\Pi} = \arg \min_{\mathbf{r} \in \mathbb{R}^q} \|\mathbf{J}\mu - \Phi\mathbf{r}\|_{\xi}^2 + \lambda \|\mathbf{r}\|_{\xi}^2. \quad (31)$$

Hence, \mathbf{r}_{Π} can be solved as

$$\mathbf{r}_{\Pi} = (\Phi^T \Xi \Phi + \lambda \mathbf{I})^{-1} \Phi^T \Xi \mathbf{J}\mu, \quad (32)$$

where the first term turns out to be nonsingular, and the existence of the inversion is guaranteed.

Applying the LSPE method along with the L_2 regulation, the synthesized algorithm for making motion decisions in the quadrotor game is listed in Algorithm 2, where $C(\cdot, \cdot)$ is the combination of operators.

4 Simulation and flight test results

In this study, six scenarios in the simulation are selected to demonstrate the efficacy of the proposed methods. The presented formulation of the one-on-one quadrotor robot game was tested in the simulation. The state transition function considering boundaries of the game field and static obstacles, and the scoring functions consisting of orientation, projected range, and terrain are examined in Scenarios 1–4. Then, the proposed motion planning algorithm is validated in Scenarios 5 and 6. After that, flight experiments are further conducted to confirm the real-time performance of the proposed decision strategy.

4.1 Scenario 1

Direct action commands are first provided to validate the state transition function (Algorithm 1). The blue quadrotor starts at the initial point

Algorithm 2 Motion planning algorithm

```

1:  $\bar{S} = \emptyset$  // The set of visited states
2:  $\mathbf{r}_0 = O_{(1+2q+C(q,2)) \times 1}$  // Initial fit coefficients
3: for each step  $k = 1, 2, \dots$  do
4:    $\mathbf{s}_k = [x_b, y_b, z_b, u_b, v_b, w_b, \psi_b, x_r, y_r, z_r, u_r, v_r, w_r, \psi_r]$ 
   // Current game state
5:    $\bar{S} = \bar{S} \cup \mathbf{s}_k$  // Update the set of visited states
6:    $\bar{\phi}(\mathbf{s}_k) = [\bar{\phi}_1(\mathbf{s}_k), \bar{\phi}_2(\mathbf{s}_k), \dots, \bar{\phi}_q(\mathbf{s}_k)]$ 
   // Feature elements
7:    $\phi(\mathbf{s}_k) = [1, \bar{\phi}(\mathbf{s}_k), \bar{\phi}^2(\mathbf{s}_k), \bar{\phi}_l(\mathbf{s}_k)\bar{\phi}_m(\mathbf{s}_k)]^T, l \neq m$ 
   // Augmented feature vector
8:   if  $\mathbf{s}_k \notin \bar{S}$  then
9:      $\mu(\mathbf{s}_k) = \mu_0(\mathbf{s}_k)$  // Initial policy
10:  else
11:     $\mu(\mathbf{s}_k) = \mu(\mathbf{s}_i = \mathbf{s}_k | \mathbf{s}_i \in \bar{S})$ 
12:  end if
13:   $\mathbf{s}_{k+1} = f(\mathbf{s}_k, \mu(\mathbf{s}_k), \omega_k)$  // Update game state
14:   $\mathbf{r}_{k+1} = \left( \sum_{t=0}^{k-1} \phi(\mathbf{s}_t)\phi^T(\mathbf{s}_t) + \phi(\mathbf{s}_k)\phi^T(\mathbf{s}_k) + \lambda \mathbf{I} \right)^{-1} \cdot \left( \sum_{t=0}^{k-1} \phi(\mathbf{s}_t)\gamma\phi^T(\mathbf{s}_{t+1}) + \phi(\mathbf{s}_k)\gamma\phi^T(\mathbf{s}_{k+1}) \right) \mathbf{r}_k$ 
   +  $\sum_{t=0}^{k-1} \phi(\mathbf{s}_t)g(\mathbf{s}_t, \mathbf{s}_{t+1}) + \phi(\mathbf{s}_k)g(\mathbf{s}_k, \mathbf{s}_{k+1})$ 
   // Update fit coefficients by the LSPE with  $L_2$ 
   // regulation
15:  if  $k$  is sufficiently large then
16:     $\tilde{J}\mu(\mathbf{s}_i) = \phi^T(\mathbf{s}_i)\mathbf{r}_k$  // Policy evaluation
17:     $\mu(\mathbf{s}_i | \mathbf{s}_i = \mathbf{s}_k) = \arg \min_{a \in A} \sum_{j=1}^n p_{ij}(a)(g(\mathbf{s}_i, \mathbf{s}_j) + \gamma \tilde{J}\mu(\mathbf{s}_j, \mathbf{r}_k))$  // Policy improvement
18:  end if
20: end for

```

$(-1.2, 0, 0.4)$ m and approaches the obstacle, while the red quadrotor flies towards the obstacle from the initial position $(-1.2, 0.7, 0.8)$ m. Both quadrotors maintain the invariant flight height. Hard constraints on translational motions are adopted; thus, the quadrotors are expected to neither exceed the scope of the game field, nor impact the obstacle. Fig. 4 illustrates that the quadrotor detours near the obstacle if the flight height is below the safety height (blue), and flies across the obstacle if the flight height is beyond the safety height (red). Both quadrotors stop at the boundaries of the game field.

4.2 Scenario 2

To verify the relationship of the S_A with the orientation angles and projected range, both quadrotors make a complete round pure yaw motion one after the other. Fig. 5 shows that the blue player achieves advantaged status with a lower value of S_A , while

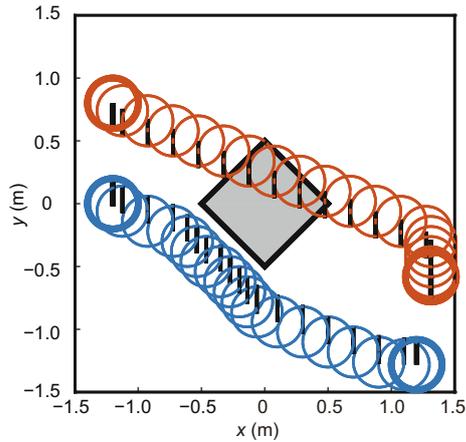


Fig. 4 Simulation results of the state transition function

Quadrotors are represented by two circles distinguished by blue and red. Solid line, starting from the center of the circle, points out the facing direction of the quadrotor. References to color refer to the online version of this figure

the red player prefers a higher value of S_A . After that, the blue quadrotor stays at the original point, whereas the red quadrotor is commanded to move perpendicular to the facing direction from two starting points with projected ranges R_1 and R_2 . Fig. 6 illustrates that it becomes insensitive to change S_A by translational motions when the distance between the quadrotors increases. Hence, it is necessary to introduce the contribution of range to encourage the quadrotors to engage in the game.

4.3 Scenario 3

To reveal the variance of S_R according to the projected range R , the red quadrotor hovers at a fixed point, while the blue quadrotor flies around the red quadrotor at different distances. Fig. 7 shows the contour diagram of the S_R , which indicates that the blue quadrotor is attracted by the red quadrotor if the projected range R is larger than the desired range R_d , whereas if R is smaller than R_d , the blue quadrotor will be pushed away from the red one.

4.4 Scenario 4

To depict the relationship of the S_H and flight height, and the influence of the obstacle, the red quadrotor hovers at a fixed point, while the blue quadrotor is commanded to lift from the ground to the ceiling of the game field outside and inside the obstacle area. Fig. 8 illustrates that the cost function reaches its minimum when the blue and red quadro-

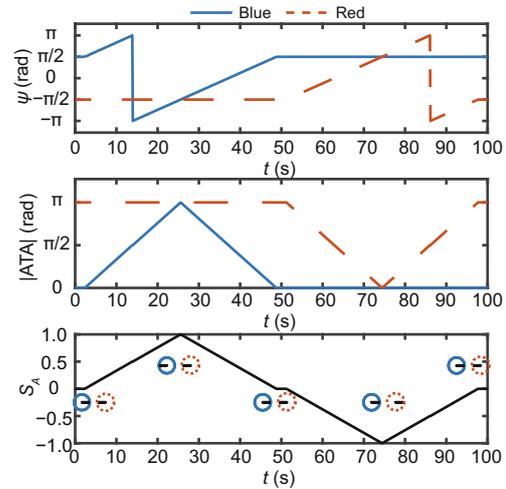


Fig. 5 The relationship of the S_A with orientation angles

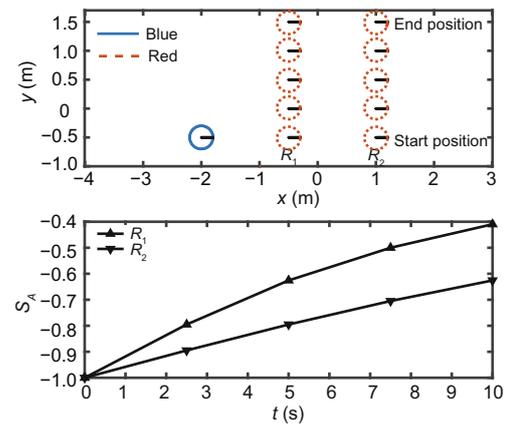


Fig. 6 Influence of the projected range to S_A

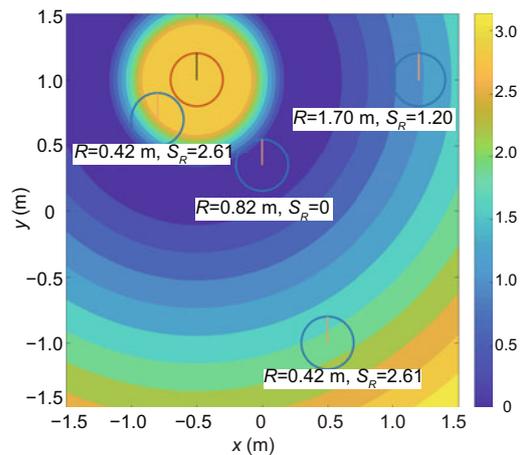


Fig. 7 The contour diagram of S_R

tors are at the same height outside the obstacle area, whereas the cost function decreases to 1.20 when the

blue quadrotor is ascending inside the obstacle area.

According to the above simulation results, the presented formulation of the quadrotor robot game is qualified to further testify the performance of the proposed decision algorithm in Scenarios 5 and 6.

4.5 Scenario 5

A simplified quadrotor game setup is tested to verify the convergence and fundamental performance of the motion planning algorithm (Algorithm 2). At the beginning, the deployed quadrotors hover at the initial positions as shown in Fig. 9. The red quadrotor remains at the start position during the simulation, while the blue quadrotor begins with an initial policy $\mu(s_i) = \psi^+$, $i = 1, 2, \dots, n$. The feature elements are chosen as $\vec{\phi} = [1 + \psi_b, 1 + ATA_b, 1 + S_A]$ for demonstration. Hence, the augmented feature vector is expressed as $\phi = [1, 1 + \psi_b, 1 + ATA_b, 1 + S_A, (1 + \psi_b)^2, (1 + ATA_b)^2, (1 + S_A)^2, (1 + \psi_b)(1 + ATA_b), (1 + \psi_b)(1 + S_A), (1 + ATA_b)(1 + S_A)]$. At the early stage of the simulation, γ is set close to zero to suppress the influence of the history value of r_k , because the number of simulation samples is largely insufficient to represent the distribution of game states by the frequencies of the visited states. As the number of samples increases, γ gradually approaches one after the distribution tends to stabilize; therefore, the long-term utility value can be learned during the policy evaluation. Fig. 10 illustrates the convergence of the fit coefficient vector r , and Fig. 11 shows the convergence of the probability vector of states ξ . Fig. 9 demonstrates the performance of the improved policy compared with the initial policy in four cases.

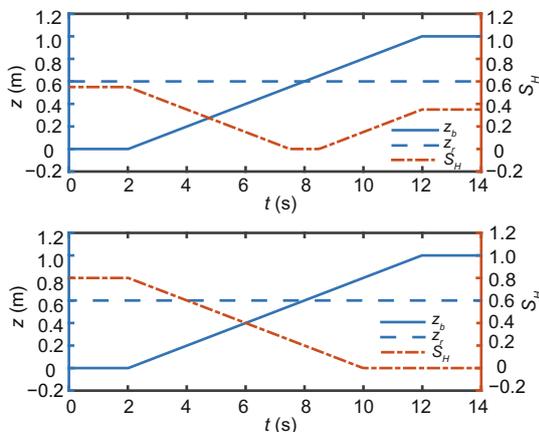


Fig. 8 The relationship of S_H and flight height and the influence of the obstacle

Initial facing directions in the first three cases are chosen as $0, -\pi/2, \text{ and } \pi$, respectively. In the fourth case, a random angle is selected as the initial condition. The results validate that the motion policy generated by the proposed algorithm can drive the quadrotor to advantageous states despite the initial positions.

4.6 Scenario 6

A synthesized scenario is finally tested to verify the performance of the proposed motion planning algorithm in a complex environment framed as the presented game formulation defined in Section 2. The blue player starts from the initial position $(-1.0, 1.0, 0.4)$ m with the facing direction of $\pi/2$, while the red player commences at the initial position $(1.0, -1.0, 0.4)$ m with the same facing direction. The initial policy of the blue player is set as

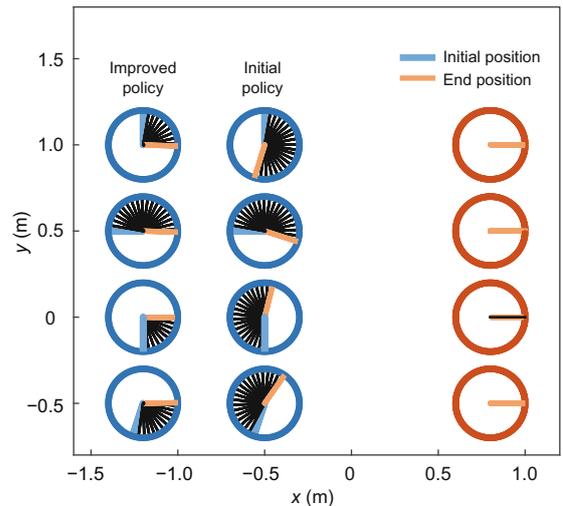


Fig. 9 Fundamental performance of the proposed decision algorithm

References to color refer to the online version of this figure

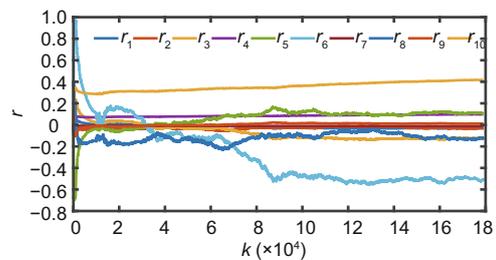


Fig. 10 Convergence process of the fit coefficient vector r

References to color refer to the online version of this figure

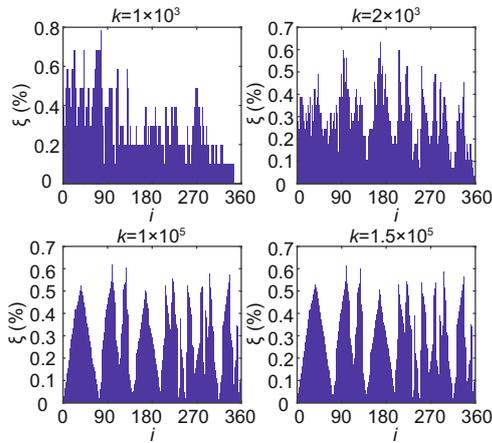


Fig. 11 Convergence process of the probability vector of states

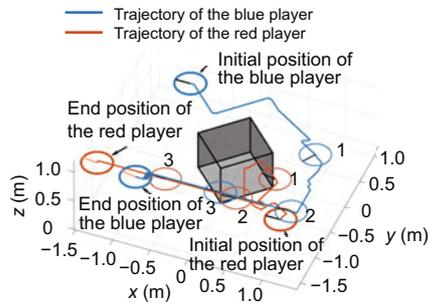


Fig. 12 Simulation trajectory of the quadrotor robot game

References to color refer to the online version of this figure

a fixed policy $\mu_0(s_k) = \psi^-$. After sufficient iterations, the improved policy generated by Algorithm 2 is adopted by the blue player to compete with the red player driven by the *-minimax algorithm. Fig. 12 presents one trajectory of the quadrotor robot game. The blue player starts at an adverse state since it is in front of the red player with its back exposed to the opponent. However, as the game proceeds, the blue player gains advantageous status over the red player. To better evaluate the performance of the algorithm, extensive simulations under the same setup are conducted, where the blue player equips with both the proposed and the *-minimax algorithms to play against the red player. The simulation results suggest that compared with the reference algorithm, the proposed algorithm can perform better and more steadily (Fig. 13).

Finally, real flight tests are conducted in the Networked Autonomous Vehicles (NAV) Lab at the Concordia University with two Quanser QDrone UAVs to further validate the real-

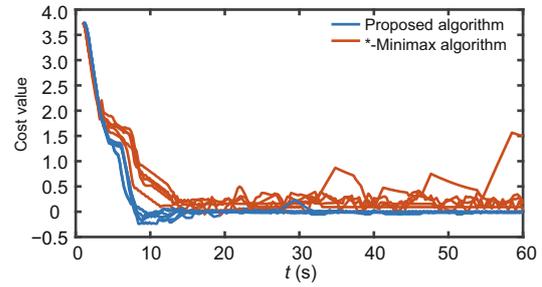


Fig. 13 The performance of the proposed algorithm in comparison to the *-minimax algorithm

References to color refer to the online version of this figure

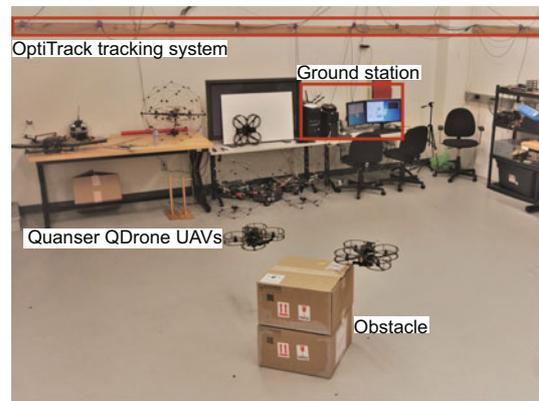


Fig. 14 Flight test conducted in the Networked Autonomous Vehicles (NAV) Lab at the Concordia University

time performance of the proposed decision strategy. Videos of the experiments are available at <https://youtu.be/Omzab2Rt2Ns>. Fig. 14 shows the scene of the flight test, which consists of a tracking system, a ground station, an obstacle, and two quadrotor UAVs. Fig. 15 shows the flight trajectory of the quadrotors in one game. The involved quadrotors start at two corners of the game field, separated by the obstacle on the ground. The commands for both UAVs are produced according to the real-time positions and motion states of the quadrotors. Therefore, the motion planning algorithm must be adequately efficient to generate rational decisions; otherwise, delayed commands may cause unexpected reactions to the varying game situations or collisions with the obstacle or the opponent UAV during the game. In the flight test, the blue player employs the proposed motion planning algorithm, while the red player uses the *-minimax algorithm as in the simulation. However, the real motion in the flight test differs from the output of the state transition model in the simulation because of the intense disturbances

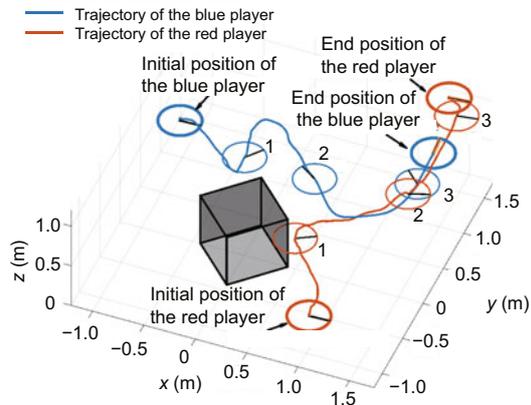


Fig. 15 Flight trajectory of the quadrotor robot game

References to color refer to the online version of this figure

caused by the rolling propellers. It can be observed in Fig. 15 that the blue player can obtain advantaged status in the game in spite of the existence of strong disturbances. The flight tests emphasize the advantage of the proposed algorithm that the generated policy is improved by the history data of the game. As a result, the mismatch of the state transition model of the UAV can be automatically corrected during the flight. Therefore, this dynamic policy can achieve a better performance compared with a fixed policy.

5 Conclusions and future work

The formulation of a velocity varying one-on-one quadrotor robot game motion planning problem has been proposed in this study. Under this formulation, a new simulation-based quadrotor motion planning algorithm using the projected policy iteration method has been proposed, and its effectiveness has been verified by extensive simulation and experimental tests. In future work, a feature refining method is to be developed in addition to the employed L_2 regulation method. Based on this, an upgraded motion planning algorithm is then to be formed to improve the capability of making rational decisions in a more sophisticated environment.

References

- Ballard BW, 1983. The *-minimax search procedure for trees containing chance nodes. *Artif Intell*, 21(3):327-350. [https://doi.org/10.1016/S0004-3702\(83\)80015-0](https://doi.org/10.1016/S0004-3702(83)80015-0)
- Bellman R, 1952. On the theory of dynamic programming. *Proc Nat Acad Sci*, 38(8):716-719. <https://doi.org/10.1073/pnas.38.8.716>
- Bertsekas DP, 1971. *Dynamic Programming and Optimal Control*. Athena Scientific, Belmont, Massachusetts, USA.
- Bertsekas DP, 2007. *Dynamic Programming and Optimal Control (3rd Ed.)*. Athena Scientific, Belmont, Massachusetts, USA.
- Bertsekas DP, 2011. Temporal difference methods for general projected equations. *IEEE Trans Autom Contr*, 56(9):2128-2139. <https://doi.org/10.1109/TAC.2011.2115290>
- Bertsekas DP, 2012. Incremental gradient, subgradient, and proximal methods for convex optimization: a survey. In: Suvrit Sra SN, Wright SJ (Eds.), *Optimization for Machine Learning*. MIT Press, Massachusetts, USA.
- Bertsekas DP, 2015. Lambda-policy iteration: a review and a new implementation. <https://arxiv.org/abs/1507.01029>
- Bertsekas DP, Tsitsiklis JN, 2000. Gradient convergence in gradient methods with errors. *SIAM J Optim*, 10(3):627-642. <https://doi.org/10.1137/S1052623497331063>
- Buşoniu L, Ernst D, de Schutter B, et al., 2010. Online least-squares policy iteration for reinforcement learning control. *Proc American Control Conf*, p.486-491. <https://doi.org/10.1109/ACC.2010.5530856>
- Efroni Y, Dalal G, Scherrer B, et al., 2018a. Beyond the one step greedy approach in reinforcement learning. <https://arxiv.org/abs/1802.03654>
- Efroni Y, Dalal G, Scherrer B, et al., 2018b. Multiple-step greedy policies in online and approximate reinforcement learning. <https://arxiv.org/abs/1805.07956>
- Fang J, Zhang LM, Fang W, et al., 2016. Approximate dynamic programming for CGF air combat maneuvering decision. *2nd IEEE Int Conf on Computer and Communications*, p.1386-1390. <https://doi.org/10.1109/CompComm.2016.7924931>
- Ghamry KA, Dong YQ, Kamel MA, et al., 2016. Real-time autonomous take-off, tracking and landing of UAV on a moving UGV platform. *24th Mediterranean Conf on Control and Automation*, p.1236-1241. <https://doi.org/10.1109/MED.2016.7535886>
- Hastie T, Tibshirani R, Friedman J, 2001. *The Elements of Statistical Learning*. Springer, New York, USA.
- Hauk T, Buro M, Schaeffer J, 2004. Rediscovering *-minimax search. *Int Conf on Computers and Games*, p.35-50. https://doi.org/10.1007/11674399_3
- Liu ZX, Zhang YM, Yu X, et al., 2016. Unmanned surface vehicles: an overview of developments and challenges. *Ann Rev Contr*, 41:71-93. <https://doi.org/10.1016/j.arcontrol.2016.04.018>
- Ma YF, Ma XL, Song X, 2014. A case study on air combat decision using approximated dynamic programming. *Math Probl Eng*, 2014:183401. <https://doi.org/10.1155/2014/183401>
- McGrew JS, 2008. *Real-Time Maneuvering Decisions for Autonomous Air Combat*. MS Thesis, Massachusetts Institute of Technology, Massachusetts, USA.
- McGrew JS, How JP, Williams B, et al., 2010. Air-combat strategy using approximate dynamic programming. *J Guid Contr Dynam*, 33(5):1641-1654. <https://doi.org/10.2514/1.46815>
- Powell WB, 2007. *Approximate Dynamic Programming: Solving the Curses of Dimensionality*. John Wiley & Sons, New Jersey, USA.

- Russell SJ, Norvig P, 2010. Artificial Intelligence: a Modern Approach (3rd Ed.). Prentice Hall, New Jersey, USA.
- Sharifi F, Chamseddine A, Mahboubi H, et al., 2016. A distributed deployment strategy for a network of cooperative autonomous vehicles. *IEEE Trans Contr Syst Technol*, 23(2):737-745.
<https://doi.org/10.1109/TCST.2014.2341658>
- Sutton RS, Barto AG, 1998. Reinforcement Learning: an Introduction. MIT Press, Massachusetts, USA.
- Thiery C, Scherrer B, 2010. Least-squares λ policy iteration: bias-variance trade-off in control problems. Proc 27th Int Conf on Machine Learning, p.1071-1078.
- Wang B, Zhang YM, 2018. An adaptive fault-tolerant sliding mode control allocation scheme for multirotor helicopter subject to simultaneous actuator faults. *IEEE Trans Ind Electron*, 65(5):4227-4236.
<https://doi.org/10.1109/TIE.2017.2772153>
- Wang B, Yu X, Mu LX, et al., 2019. Disturbance observer-based adaptive fault-tolerant control for a quadrotor helicopter subject to parametric uncertainties and external disturbances. *Mech Syst Signal Process*, 120:727-743.
<https://doi.org/10.1016/j.ymsp.2018.11.001>
- Yu HZ, 2010. Convergence of least squares temporal difference methods under general conditions. 27th Int Conf on Machine Learning, p.1207-1214.
- Yu HZ, 2012. Least squares temporal difference methods: an analysis under general conditions. *SIAM J Contr Optim*, 50(6):3310-3343.
<https://doi.org/10.1137/100807879>
- Yuan C, Zhang YM, Liu ZX, 2015. A survey on technologies for automatic forest fire monitoring, detection, and fighting using unmanned aerial vehicles and remote sensing techniques. *Can J Forest Res*, 45(7):783-792.
<https://doi.org/10.1139/cjfr-2014-0347>