

Path Planning based on Fuzzy Bayes and Deep Q-Network Algorithm

Feifei Yin¹, Dong Liu^{1, a}, Yu Gong¹

¹School of North China Electric Power University, Baoding 071000, China.

^a1206776367@yahoo.com

Abstract

Virtual assembly with the aid of computer virtual reality technology can realize the design and planning of assembly process, improve product production efficiency and reduce economic cost. Path planning is an important development direction of virtual assembly, the research of path planning technology in virtual assembly environment is very important for assembly path design in complex environment. In this paper, a deep q-network algorithm based on Fuzzy Bayes is proposed and tested in virtual assembly environment, concluded that the depth q-network algorithm based on Fuzzy Bayes has better trafficability and planning efficiency in the complex environment with narrow space.

Keywords

Path planning; deep reinforcement learning; fuzzy Bayesian decision algorithm; fuzzy Bayesian deep Q network.

1. Preface

This year, With the rapid development of virtual reality technology, it provides a new solution for industrial simulation assembly technology. Virtual reality technology is applied to the assembly design of mechanical and electrical products, and the virtual assembly technology is born. Virtual assembly is an important part of virtual manufacturing, By using computer to simulate the real assembly operation and demonstrate the assembly process with visual tools, problems can be found as soon as possible and a more reasonable assembly scheme can be formulated. This method saves time and effort, improves performance, easily finds problems, and can be modified in time. It greatly shortens the R & D and manufacturing cycle of products, reduces the cost of products and improves the competitiveness of products. In virtual assembly, assembly path planning technology can realize the automatic generation of path, which can greatly reduce the economic cost. Therefore, the research of automatic path planning technology has practical significance for virtual assembly. In this paper, we use reference [1] to find the shortest path of the document; In reference [2], an improved A* algorithm is proposed to solve the robot path planning problem, and the experimental results show that the success rate of the improved algorithm is higher than that of the original A* algorithm.

Reference [3] searched randomly assigned nodes in the space, and calculated and selected the shortest route. In reference [4], a new particle swarm optimization algorithm based on non-uniform Markov chain is constructed, which provides a new method for intelligent robot path planning. In reference [5], the path planning ability of robot is trained by using full convolution neural network and fast extended random tree algorithm. Reference [6] combines human-computer interaction guidance and path solving algorithm to guide the generation of assembly path. In reference [7], the artificial potential field local optimization algorithm is introduced into ant colony algorithm to improve the convergence speed of global path planning. In reference [8], an improved particle swarm optimization algorithm is proposed to optimize the path of solder

joint detection of circuit board. Throughout the research status at our country and abroad, path planning technology has made great progress, but in the complex environment, how to successfully model, how to plan the path is still a very important problem. In addition, the research on virtual assembly path planning is less successful, and the practical application is even less.

In this paper, an improved deep reinforcement learning algorithm is proposed to solve the problem of narrow space in complex environment and through the experiment to compare and verify.

2. Algorithm Model

In this paper, the deep Q-network algorithm in deep reinforcement learning is used to solve the path planning problem of virtual assembly, and the fuzzy Bayesian decision-making algorithm is proposed to solve the problem of exploration and utilization. Firstly, the prior knowledge in exploration is obtained by using fuzzy comprehensive evaluation method to synthesize various factors. Then the Bayesian decision algorithm updates the posterior distribution based on the known prior probability distribution. Finally, actions are selected according to posterior distribution to generate decisions to optimize the path planning of virtual assembly.

2.1. Deep Q Network Algorithm

Deep Q network algorithm [13] is mainly composed of Q-learning and deep convolutional neural network (CNN) [10,11,12]. Q-learning is a reinforcement learning algorithm with different strategies.

The learning process of Q-learning is as follows:

In the first stage: the agent chooses action a to interact with the environment according to the ϵ -greedy algorithm under the condition of randomly given state s ;

In the second stage: after the agent performs action a , the current state changes and enters the next state. At the same time, the environment will give the agent an immediate feedback (reward or punishment).

In the third stage, the R matrix is constructed with the state as the row and the action as the column, and the real-time feedback of the environment to the agent after the action a is executed in the state s is stored in the R matrix;

The fourth stage: according to R matrix, P matrix is calculated to guide the agent's action. The P matrix is composed of state and action key value pairs (state action value function, Q value). The greedy strategy is used to evaluate and improve the Q value, and the formula (1) is used to update the Q value and converge to the optimal Q value;

The fifth stage: according to the P matrix, select the action with the largest value function (Q^* , the best return) in each state, and generate the action sequence. The update formula of state action value function is as follows:

$$Q(s_t, a) \leftarrow Q(s_t, a) + \alpha [R(s_t, s_{t+1}, a) + \gamma \max_{a' \in A} Q(s_{t+1}, a') - Q(s_t, a)] \quad (1)$$

Where s_t and s_{t+1} represent the current state and the next state respectively, and a represents the action selected by the agent, $Q(s_t, a)$ refers to the state action Q value of agent a in state s_t , α represents learning rate, $R(s_t, s_{t+1}, a)$ refers to immediate feedback given by environment when agent performs action a and current state is transferred from s_t to s_{t+1} . γ is discount factor, which indicates the influence of future feedback on the current. The smaller γ is, the more attention the agent pays to the immediate feedback brought by the recent decision, and $\max_{a' \in A} Q(s_{t+1}, a', \theta)$ represents the optimal value function for predicting the next state.

In Q-learning, a Q-function $Q(s, a)$ is defined to represent the feedback that can be obtained by taking action a in state S . The Q-value is constantly updated by iteration. If the Q-function is accurate enough and the environment is certain, the strategy of selecting the maximum Q-value action can be adopted.

The traditional action value P matrix is mainly realized by the lookup table [14], and the Q value is stored in a Q table. However, when the state increases, the number of agents' actions increases exponentially, which makes the data difficult to store.

In practice, some unprocessed images are often used to represent the state, so it is difficult to apply Q-learning to practical problems. However, deep convolution neural network can extract feature information from images, abstract and classify them, etc. So we use deep convolution neural network to simulate Q function. Therefore, DQN algorithm uses CNN to solve the dimension disaster problem of P matrix in Q-learning. The process is as follows:

Firstly, when inputting the state, the current state is represented by the images of different observable environments and their position information of the agent. Then Use CNN to extract abstract features for each different state, so as to reduce the dimension of the original high-dimensional state space and make it easy to store; Secondly, CNN evaluates and selects the strategy of agents every k frames. Among them, the skip frame keeps the original selection strategy unchanged; Then, the exploration environment stores the collected data in the form of memory unit $(s_t, a_t, r_{t+1}, s_{t+1})$ in the experience playback cache. In this algorithm, s_t and a_t are the States and selected actions in time step 't', and then randomly select samples from the experience playback buffer to train the neural network; Finally, the CNN with parameter θ_i is used to approximate the Q value, and the approximation value function $f(s, a; \theta_i)$ is established to approximate the optimal value function Q^* , where 's' and 'a' are the current state and action of the agent respectively, and θ_i is the parameter of the i-th iteration, and the corresponding Q value of each state is output.

In DQN algorithm, the mean square error is used to define the target loss function, and the target Q value in Q learning is taken as the label. The target Q value is represented by $R(s_t, s_{t+1}, a) + \gamma \max_{a \in A} Q(s_t, a)$ in formula (1). The target loss function value is calculated according to the deviation between the target Q value and the predicted output. The loss function is optimized by the synchronous training method of Q-learning algorithm and random gradient descent method to minimize the loss function value. The loss function is defined as follows:

$$L(\theta) = E[(r + \gamma \max_{a'} Q(s_{t+1}, a', \theta) - Q(s_t, a, \theta))^2] \quad (2)$$

Among them, s_t and s_{t+1} represent the current state and the next state respectively, a and a' respectively represent the action selected by the current state and the action selected by the next state, θ is the parameter of CNN, 'r' is the instant feedback, $\gamma (0 < \gamma < 1)$ is the discount factor, $\max_{a'} Q(s_{t+1}, a', \theta)$ is the optimal value function of the prediction after executing the action a' on the next state, and $Q(s_t, a, \theta)$ is the Q value of the current prediction output.

2.2. Fuzzy Bayesian Decision Algorithm with Prior Knowledge

Quantitative evaluation of various factors to make decisions. Fuzzy Bayesian decision algorithm with prior knowledge [15,16,17], The steps are as follows:

Step 1: construct the factor set U of fuzzy comprehensive decision-making method, and the factor set u is the set of all factors affecting the decision-making. Set the action set a and state set s of Q learning, and divide the state set s according to factor set U ;

Step 2: construct fuzzy evaluation matrix E_f and weight set W according to experts' experience. According to the factor set U and fuzzy evaluation matrix E_f , through the weight set W , all States $s_i, i=1,2,\dots,k$. The superiority degree of each decision under K is given by the superiority vector $B_i, i=1,2,\dots,K$ means. Then, the B_i of each state s_i is normalized, and the result is used as the prior knowledge of Q learning to initialize the Q value of state s_i ;

Step 3: start Q learning. In the time step t , according to the current state s_t , select action a_j to reach the new state s_{t+1} , get an immediate return $r(s_t, a_j)$, and update the Q value;

In the learning process, Boltzmann method is used to calculate the probability $p(a_j)$ to select the action randomly. Probability $p(a_j)$ is defined as follows:

$$p(a_j) = \frac{\exp\left[\frac{Q(s, a_j)}{T}\right]}{\sum_k \exp\left[\frac{Q(s, a_k)}{T}\right]} \quad (3)$$

Where $Q(s, a_j)$ is the Q value of the state s -action a_j pair, T is the temperature parameter in the annealing process. The larger T is, the greater the probability of random sampling is;

Step 4: the posterior probability distribution is calculated by Bayes formula, and the improved Bayes formula (5) is obtained by combining the action and state in Q -learning with formula (4), The definition is as follows:

$$\theta_{t+1} = \theta_t + \alpha \left[r + \gamma \max_{a'} Q(s', a'; \theta) - Q(s, a; \theta) \right] \nabla Q(s, a; \theta) \quad (4)$$

$$p(a_j | s_t) = \frac{p(s_t | a_j) \cdot p(a_j)}{p(s_t)} \quad (5)$$

Among them, $p(s_t)$ is the probability of occurrence of state s_t , $p(s_t | a_j)$ is the probability of selecting action a_j , $p(s_t | a_j)$ is the probability of state transferring to s_t after selecting action a_j , and $p(a_j | s_t)$ is the probability of agent selecting action a_j in state s_t ;

Step 5: compare the Bayesian risk of making a decision immediately in the state s_t with the expected posterior Bayesian risk from the next state observation. The loss function is calculated by formula (2), and the minimum posterior risk is obtained by optimizing the loss function according to the synchronous training method of Q -learning algorithm and random gradient descent method;

Step 6 repeat the above steps until the end of the study. The optimal action sequence is selected to generate the decision.

3. Analysis of Experimental Results

The experiment is based on OpenAI Gym environment library and implemented by TensorFlow, a simulation environment with unreachable state is established randomly based on OpenAI Gym environment library.

The experiment shows that although the fuzzy Bayesian deep Q network algorithm takes more time to explore the environment and consumes a certain amount of time to explore the

environment, the path will be re optimized immediately. When the number of steps reaches 120, the path tends to be stable.

In order to verify the planning efficiency of deep q-network in virtual assembly path planning in narrow space, this paper compares the experimental results of the traditional fast expanding random tree algorithm(Rapid-exploration Random Tree, RRT) [18,19,20] and the improved fuzzy Bayesian depth q-network in the same simulation environment.

Table 1. Comparison of planning efficiency between RRT algorithm and fuzzy Bayesian deep Q network algorithm

	50%	40%	30%	20%	10%
RRT[25]	40.0184527	39.4535781	34.5006577	31.5828547	321.0748215
RRT[100]	3.1327235	3.0223597	3.2958478	3.9965231	30.6325810
FB-DQN	2.6912584	2.7098782	2.8379538	2.6258714	2.5997448

The above table shows the efficiency comparison chart of the two algorithms. It can be seen from the comparison that the planning efficiency of RRT algorithm depends on the setting of search step size. However, in this experiment, when the step size is too large, it will deviate from the minimum width of the free space in the global map, resulting in planning failure. Although the path planning time is reduced to a certain extent, the success rate also decreases. However, the fuzzy Bayesian depth Q-network algorithm does not need to set the search step size, and has stable planning efficiency and short planning time under different global traffic degrees. In the narrow space, using the fuzzy Bayesian depth Q network algorithm for path planning has better planning efficiency. Therefore, in the narrow space, using deep Q-network algorithm to solve the path planning problem of virtual assembly has better trafficability and planning efficiency.

4. Epilogue

In this paper, fuzzy Bayesian decision-making algorithm and fuzzy Bayesian depth Q-network algorithm with prior knowledge are proposed. The experimental results show that the proposed algorithm has good trafficability and planning efficiency for virtual assembly path planning, and achieves the expected goal.

Acknowledgements

The authors acknowledge the Fundamental Research Funds for the Central Universities (Grant:2018 MS078)

References

- [1] Singh Y, Sharma S, Sutton R, et al. Optimal Path Planning of an Unmanned Surface Vehicle in a Real-Time Marine Environment using a Dijkstra Algorithm[C]// The International Conference on Marine Navigation and Safety of Sea Transportation. 2018:399-402.
- [2] Fu B, Chen L, Zhou Y, et al. An improved A* algorithm for the industrial robot path planning with high success rate and short length[J]. Robotics & Autonomous Systems, 2018.
- [3] Altinoz O T, Yanar T A, Ozguven C, et al. Improved non-Probabilistic Roadmap method for determination of shortest nautical navigation path[C]// International Conference on Electrical and Electronic Engineering. IEEE, 2017:261-266.
- [4] Zeng N, Zhang H, Chen Y, et al. Path planning for intelligent robot based on switching local evolutionary PSO algorithm[J]. Assembly Automation, 2016, 36(2):120-126.

- [5] Pérezhigueras N, Caballero F, Merino L. Learning Human-Aware Path Planning with Fully Convolutional Networks[C]// IEEE International Conference on Robotics and Automation. IEEE, 2018.
- [6] Flavigne D, Taix M, Ferre E. Interactive motion planning for assembly tasks[C]// The, IEEE International Symposium on Robot and Human Interactive Journal of agricultural machinery Communication, 2009. Ro-Man. IEEE, 2009:430-435.
- [7] Liu Jianhua, Yang Jianguo, Liu Huaping, Geng Peng, Gao Meng. Global path planning method for mobile robot based on potential field ant colony algorithm[J]. Journal of agricultural machinery, 2015, 46(09):18-27.
- [8] Lillicrap T P, Hunt J J, Pritzel A, et al. Continuous control with deep reinforcement learning[J]. Computer Science, 2015, 8(6):A187.
- [9] Wang Z, Schaul T, Hessel M, et al. Dueling network architectures for deep reinforcement learning[J]. 2015.
- [10] Elbaz G, Avraham T, Fischer A. 3D Point Cloud Registration for Localization Using a Deep Neural Network Auto-Encoder[C]// Computer Vision & Pattern Recognition. 2017.
- [11] He K, Zhang X, Ren S, et al. Deep Residual Learning for Image Recognition[C]// 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). 2016:770-778.
- [12] Llinas J., Hall D.L. An introduction to multisensor data fusion[J]. 1997, 85(1).
- [13] Pan Y, Zhang L, Li Z W, et al. Improved Fuzzy Bayesian Network-Based Risk Analysis With Interval-Valued Fuzzy Sets and D-S Evidence Theory[J]. IEEE Transactions on Fuzzy Systems, 2019, PP(99):1-1.
- [14] RODRIGUEZ S, TANG X, LIEN J M, et al. An obstacle-based rapidly-exploring random tree: IEEE International Conference on Robotics and Automation[C].
- [15] Zhuxia, Chenrenwen, Xudongxia, etc. Solder joint inspection path planning method based on Improved Particle Swarm Optimization[J]. Journal of instrumentation, 2014(11):2484-2493.
- [16] Shuokaiping, Zhangxiaoshun, Xutao, etc. Joint optimal scheduling of multi energy systems based on knowledge transfer Q-learning algorithm[J]. Power system automation, 2017, 41(15):18-25.
- [17] Wang Xuesong, Zhu meiqiang, Cheng Yuhu. Reinforcement learning principle and its application [M]. 2014.
- [18] Cheng Xiaoping. Fuzzy Bayesian networks based on α - cut sets [J]. Computer science, 1999, 000 (006): 65-66,90.
- [19] Liu Xiaoqian, Zhang Hui, Wang Yingjian. Path planning algorithm based on improved RRT [J]. Automation technology and application, 2019 (5): 96-100.
- [20] Liu Chengju, Han Junqiang, Ankang. Dynamic path planning of robocub robot based on improved RRT algorithm [J]. Robot, 2017 (39): 8-15.