



# A Computer Game Decision Interpretation Method Based on Salient Features

Haodong Feng and Shuqin Li<sup>(✉)</sup>

Beijing Information Science and Technology University, Beijing, China  
Lishuqin\_de@126.com

**Abstract.** Artificial intelligence has greatly improved the efficiency of industry and life, but its algorithmic framework is often a black box, leading to a lack of understanding of how computers give decision results. So, the interpretability of artificial intelligence has been widely concerned in recent years. In this paper, decision interpretation research is done in the field of computer games to try to make humans understand the decision results of game agents. In this paper, three indicators are designed to analyse the game features one by one, which are threatening, relevance and specificity. Finally, all the salient features associated with the decision are given to explain the agent's decisions. Experiments are conducted on two computer games, Surakarta and Mahjong. It is found that salient features not only help humans to understand the behaviours of the agents, but also speed up the decision making of human players in the games. This shows that this paper achieves interpretability of decision making through salient features in the field of computer games.

**Keywords:** Explainable Artificial Intelligence · Computer Game · Decision Explanation

## 1 Introduction

There are problems of inherent algorithmic black boxes and opaque system information in artificial intelligence decision-making process, which can only give the final decision but not an explanation of the decision behaviours. This leads to correct but incomprehensible results of AI algorithms, which hinders the further development of AI [1][4]. Interpretability of AI refers to explaining the techniques and methods used to build AI applications, whereby people in order to understand the reasons why they make particular decisions. Interpretability of AI will help AI products to be landed, enhance public confidence and improve management. Therefore, the interpretability of AI models has become a topic of greater public interest in recent years [9]. Computer games are known as the “fruit flies” in the field of artificial intelligence, and various emerging technologies can be experimented on computer games. Therefore, in this paper, the computer game is used as a vehicle for the explanatory study of the decision making of an intelligent body.

Computer games can be divided into complete information games and incomplete information games according to the degree of information knowledge in the game process. Currently, there have been attempts to explain the behaviours of agents in the field of complete information games. In 2018, Greydanus et al. derived salient diagrams to explain the behaviours of reinforcement learning agents by applying Gaussian fuzzy to different parts of the input image [3]. The core of the algorithm is to compute the difference between the value function and the policy vector between the original and the perturbed state. The method has achieved good results on Atari games. In the same year, Iyer et al. used the difference in action values between the original and perturbed states to derive the feature salient map [5], which gives a decision interpretation for chess. In 2019, Puri et al. proposed a perturbation-based method for generating black-box smart body salient maps, SARFA [8], which uses exclusive and relevant feature attributes to interpret smart body actions. In addition, the salient map generated by SARFA can also provide valuable hint information when humans play chess and help them quickly choose the correct move when solving chess puzzles. 2021, Liu et al. [7] obtained a more comprehensive salient map by performing saliency analysis on the features of blank regions.

It can be found that the mainstream research objects of the current academic community are complete information games. In this paper, we will try to design a generalized computer game decision interpretation model so that it is applicable to most computer game systems. This model will be used to explain the decisions in the complete information game “Surakarta” and the incomplete information game “Mahjong”. It is found that the salient features given by the model explain well the decisions of computer game agents and help human players to make decisions more quickly in the game.

## 2 Related Work

### 2.1 Common Computer Game Decision Methods

As early as the beginning of the 21st century, scientific institutions and scholars have proposed the Minimax algorithm, the  $\alpha$ - $\beta$  pruning algorithm, and the upper confidence bound apply to tree with the situation evaluation algorithm to solve the computer game problem. For the characteristics of different computer games, scholars proposed different ways to perform the evaluation. After that, artificial neural network techniques such as deep reinforcement learning became more and more mature, and researchers began to try to apply them to the field of computer games. Computer game agents based on deep reinforcement learning usually use deep neural networks to fit the Q-value function and select the final decision by the Q-value function. It can be found that the core of common computer game decision methods is to evaluate the winning probability of the game participants and select the action with the highest value of the situation after the decision. This paper will focus on this decision concept for explanatory work.

## 2.2 Surakarta

Surakarta is a two-player complete information game and is a popular research object in the field of complete information games. The board consists of a  $6 \times 6$  square grid with 8 arcs on the corners. To capture the opponent's pieces, they must pass through an arc tangent to the path. When all the pieces on one side are captured, the side with the remaining pieces wins. In this paper, we specify that the horizontal coordinates of the pieces are A to F and the vertical coordinates are 1 to 6.

## 2.3 Mahjong

As one of the national treasures of Chinese culture, mahjong has a rich charm and connotation, and is characterized by its long-standing oriental culture. The object of this paper is popular mahjong, which is simplified based on the Chinese national standard mahjong. It is characterized by a large number of tile combinations and many types, covering almost all mahjong types, which is a great test of mahjong skills and strategies. There are three types of tiles in the deck: Character, Bamboo and Dot, with points from 1 to 9, each with 4 tiles, for a total of 108 tiles. The dealer starts with 14 tiles and the other three start with 13 tiles. In the course of the line, you can Chow, Pong, Kong and Waiting, and all of them can get direct benefits. The priority of the process is: Hu > Kong > Pong > Chow. When a player wins or the tiles were all taken out, the game ends and the total score is given.

# 3 Method

## 3.1 Impact of Feature Perturbations on Decision Actions

### 3.1.1 Threatening Analysis

There are three criteria for judging the salience of features for decision making actions, the first one is threatening. In computer gaming process is accompanied by offense and defense. Good players and agents will focus on features that pose a serious threat to them. For example, in Surakarta, they focus on enemy pieces that can capture our pieces, and in mahjong, they focus on tiles that may allow their opponents to win. Depending on the rules of the computer game, different threatening calculations can be developed. In this paper, the features that pose a threat to us are directly defined as salient features.

### 3.1.2 Relevance Analysis

The second criterion for the significance of a feature for a decision action is relevance. If changing a feature also has a large impact on actions other than the decision action, it indicates that the feature has low decision relevance to the decision action. Therefore, the relevance of the feature to the decision action can be judged by comparing the range of value change of non-decision actions after feature perturbation.

The correlation is calculated as shown in Eq. 1. where  $R(f)$  refers to the correlation between feature  $f$  and the decision action derived after perturbing feature  $f$ ,  $a$  refers to the executable action,  $A$  refers to the set of executable actions,  $\hat{a}$  refers to the decision

action,  $V(a)$  refers to the value function valuation or  $Q$ -value of action  $a$ , and  $V_f(a)$  refers to the value function valuation or  $Q$ -value of action  $a$  after perturbing feature  $f$ .

$$R(f) = KL\left(\frac{V(a)}{\sum_{a' \neq \hat{a}} V(a')} \parallel \frac{V_f(a)}{\sum_{a' \neq \hat{a}} V_f(a')}\right) \quad a \in A, a \neq \hat{a} \quad (1)$$

In a statistical sense, KL scatter can be used to measure the degree of difference between two distributions. We use relative entropy to measure the change in the value of non-decisional actions before and after feature perturbation, and a high value means the feature is less relevant to the decision action, while a low value means the feature is more relevant to the decision action.

### 3.1.3 Specificity Analysis

The third criterion for the significance of a feature for a decision action is specificity. If the effect of changing a feature on a decision action is much greater than the effect on other actions, then the feature has high specificity for the decision action. Therefore, the specificity of a feature for a decision action can be judged by comparing the change in value of the decision action with the change in value of the non-decision action after the feature perturbation.

The specificity is calculated as shown in Eq. 2. Where  $S(f)$  refers to the specificity of feature  $f$  with decision action after perturbing feature  $f$ ,  $n$  refers to the number of non-decision actions,  $a'$  refers to the executable actions,  $A$  refers to the set of executable actions,  $\hat{a}$  refers to the decision action, and  $V_f(a)$  refers to the valuation of value function or  $Q$  value of action  $a$  after perturbing feature  $f$ . Higher values represent higher specificity of the feature to the decision action and lower values represent lower specificity of the feature to the decision action.

$$S(f) = \frac{n \times (V_f(\hat{a}) - V(\hat{a}))}{\sum_{a' \neq \hat{a}} (V_f(a') - V(a'))} \quad a' \in A \quad (2)$$

## 3.2 Screening for Salient Features

For different computer games, the types and numbers of features are different. For example, in Surakarta, features can be classified into categories such as our pieces, opponent's pieces, position with pieces, position without pieces, etc. In popular mahjong, features can be classified into categories such as our hand, our discard, opponent's discard, etc. Therefore, the features of the target computer game need to be classified and ranked first. Then the features are analysed for threatening in turn, and if the features are threatening, they are directly defined as salient features. Otherwise, a reasonable perturbation method is chosen for the features, and the new state after perturbing the features is obtained, and relevance analysis and specificity analysis are performed on them to derive the significance degree of the features. After analysing the significance of all features, the significance features can be given to explain the agent's decision. As shown in Fig. 1.

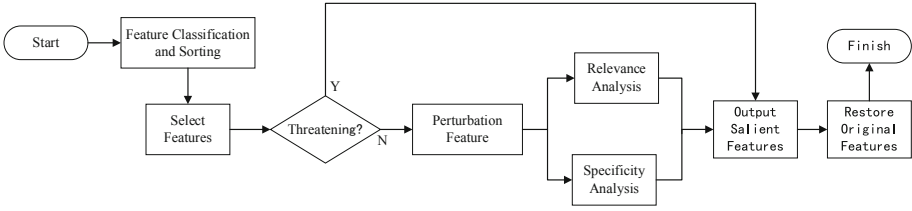


Fig. 1. Screening for salient features process.

According to the definition of KL scatter, the larger the  $R(f)$ , the lower the relevance of feature  $f$  to the decision action. And the larger the value of  $S(f)$ , the higher the specificity of feature  $f$  for the decision action. The final derived saliency calculation is shown in Eq. 3.

$$O(f) = \frac{2 \times S(f) \times \exp R(f)}{1 + S(f) \times \exp R(f)} \quad (3)$$

## 4 Experiment

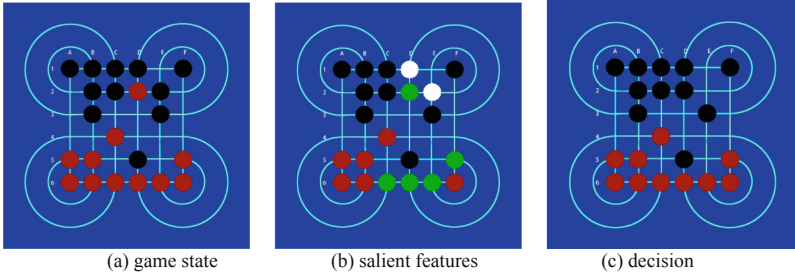
Unlike ordinary computer game studies, the explanatory work in this paper cannot be judged by comparative tests such as win rate and score to determine its merits. In this paper, three metrics are designed to verify whether the salient features given by the model are valid: comprehensiveness, correctness, and utility. In addition, in all the experiments conducted in this paper, there is a premise that only the situation value assessment function of the agents can be accessed.

### 4.1 Select Game Agents

Excellent computer gaming agents are the basis for completing the experiments. To confirm that the method in this paper has some generality, this section takes Surakarta chess in complete information game and Mahjong in incomplete information game as the experimental objects. In this paper, we refer to the articles on Surakarta chess by Li et al. [6], and construct a Surakarta gaming agent based on the Minimax Algorithm and Alpha-Beta pruning algorithm. For mahjong, the mahjong program that finished second in the 2020 Chinese Computer Game Championship [2] was chosen for this paper, and its core method is deep reinforcement learning.

### 4.2 Surakarta Agent Decision Explanation

The output of the decision interpretation salient features for the Surakarta game is shown in Fig. 2. Figure 2(a) shows the state of the game before Black's decision; Fig. 2(b) shows the salient features before Black's decision, white shows the salient features in Black and green shows the salient features in Red; Fig. 2(c) shows the state of the game after Black's decision. According to the rules of the game, the salient features of the current game state are the discs in the 1D, 2D, 2E, 5F, 6C, 6D and 6E squares. Obviously, the output of the model is correct.



**Fig. 2.** Surakarta agent decision explanation.

### 4.3 Mahjong Agent Decision Explanation

The output of decision interpretation salience features for popular mahjong is shown in Fig. 3. The tiles on the table are the discarded tiles of the corresponding players, which is a public state available to all. The rest of the information is the private state of each player. Figure 3(a) shows the state of the game before player 2's decision, and the salient features are marked by the red box; Fig. 3(b) shows the state of the game after player 2's decision. According to the rules of the game, the salient features of the current game state are 1, 2, 4, 5, 6, 7, 8 Dots, and 3, 5, 6, 8 Characters. Obviously, the output of the model is correct.

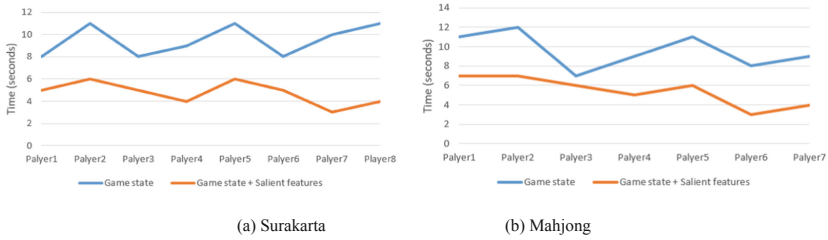
### 4.4 Match Tasks

The first validation metric in this paper is the comprehensiveness of the salient features given by the model. In each situation, human players participating in the experiment are asked to give features related to the decision making of the agents, which are matched with the salient features given by the model. The experiment was divided into 20 groups, half of which presented the salient features given by the model and the other half presented randomly selected features. If participants responded with a match for the correct salient features and a mismatch for the random salient features, it means that the salient features given by the model are comprehensive. The experimental results are shown in Table 1. This proves that the salient features given by the model are comprehensive in each game situation.

### 4.5 Explain Tasks

The second validation metric is whether the salient features given by the model are correctness. If the salient features given by the human player match the salient features given by the model, let the human player enter the explanation task. In this subsection the human player is asked to explain the agents' decisions based on the salient features given by the model. At the end of the experiment, it was found that all players could explain the agents' decisions based on the salient features, and they matched the corresponding game rules. This indicates that the salient features given by the model in the text are correct.





**Fig. 4.** The effect of salient features on the speed of decision making.

can be seen made faster decisions, which indicates that the salient features given by the model are practical.

## 5 Conclusion

To address the interpretability problem in the field of artificial intelligence, this paper makes an attempt in the field of computer games. The salient features associated with the decision making of the intelligences are given by the three main indicators of threatening, relevance, and specificity. It not only shows the decision basis of agents, but also improving the speed of human decision making. The future research can try to classify the saliency of features and refine the decision interpretation level of computer games to help humans better understand the behaviours of agents.

## References

- Adadi, A., & Berrada, M. (2018). Peeking inside the black-box: a survey on explainable artificial intelligence (XAI). *IEEE access*, 6: 52138-52160.
- Gao S. (2021). Research on mahjong machine game strategy. Beijing Information Science & Technology University.
- Greydanus, S., Koul, A., Dodge, J., & Fern, A. (2018, July). Visualizing and understanding atari agents. In *International conference on machine learning*, 1792-1801. PMLR.
- Gunning, D., Stefik, M., Choi, J., Miller, T., Stumpf, S., & Yang, G. Z. (2019). XAI—Explainable artificial intelligence. *Science Robotics*, 4(37): eaay7120.
- Iyer, R., Li, Y., Li, H., Lewis, M., Sundar, R., & Sycara, K. (2018, December). Transparency and explanation in deep reinforcement learning neural networks. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, 144-150.
- Li, S., Li, J., & Han, Y. (2012). Study on the evaluation function in the Surakarta. *Journal of Beijing Information Science & Technology University*, 27(06): 42-45.
- Liu, H., Zhang, X., & Diao, Z. (2021). An Interpretation Method of Decision Basis for the RL Agent of Chess. *Journal of Chongqing University of Technology*, 35(12): 140-146.
- Puri, N., Verma, S., Gupta, P., Kayastha, D., Deshmukh, S., Krishnamurthy, B., & Singh, S. (2019). Explain your move: Understanding agent actions using specific and relevant feature attribution. *arXiv preprint arXiv:1912.12191*.
- Tomsett, R., Braines, D., Harborne, D., Preece, A., & Chakraborty, S. (2018). Interpretable to whom? A role-based model for analyzing interpretable machine learning systems. *arXiv preprint arXiv:1806.07552*.



**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution-NonCommercial 4.0 International License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits any noncommercial use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

