

An Intentional AI for Hanabi: Reproduction

MJ Johns

Adam Khan

Introduction

In this reproduction, we consider the paper “*An intentional AI for hanabi*” (Eger, Martens, and Cordoba 2017) and propose the following research questions which reproduce and extend the findings of the prior paper:

- RQ1: Does an intentional AI agent perform better in a collaborative game with a human than a baseline AI agent? (as measured by score)
- RQ2: Is the player’s perception of the AI agent (e.g. perceived skill and/or enjoyment of playing together) correlated with higher score, regardless of implementation?
- RQ3: Can we use Machine Learning to predict human responses to improve game-playing AI for more compatible cooperation?

Background and Motivation

Game-playing AI has proven to be a useful test-bed for exploring decision-making AI implementations, and has attracted a wide range of researchers and research interests (Yin et al. 2023; Hu et al. 2024). Similarly, AI within games offer a valuable opportunity to understand human perception of interacting, and collaborating, with AI (Laird 2002; Yannakakis 2005). As we rapidly approach a state where humans and AI interact and cooperate on a daily basis, it is becoming increasingly important to (1) develop AIs that meaningfully cooperate with humans, and (2) understand what aspects of human-AI interaction are enjoyable, challenging, or unpleasant for humans. With this paper, we aim to contribute to further understanding AI implementations and human-AI interaction through an analysis of an intentional AI for the game of hanabi.

Dataset Description

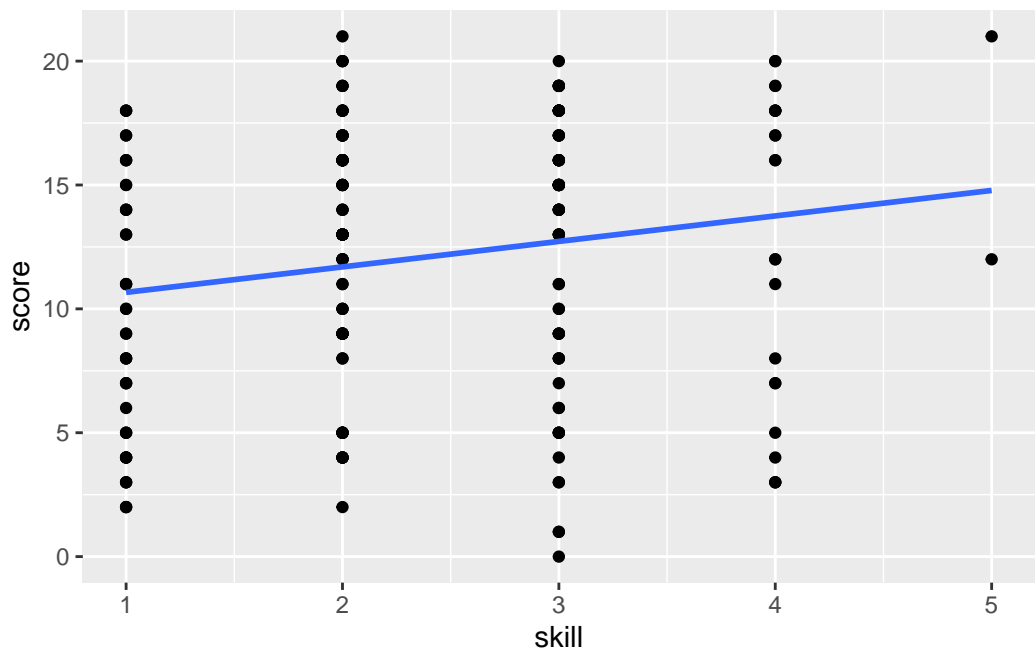
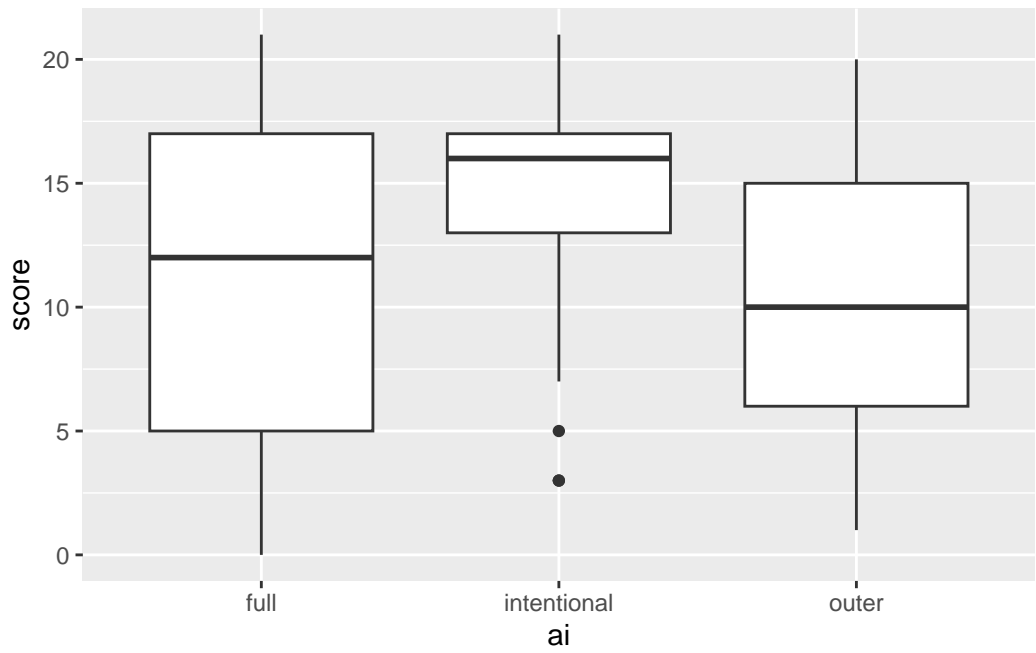
We are using two datasets for our analysis. First is the results of 240 games along with survey information from participants, which is publicly available on GitHub: <https://github.com/yawgmoth/HanabiData/> The second dataset is the turn-by-turn game logs from all games scoring over 20 points. While this data is not publicly available, we were granted access by one of the original authors. The game logs use one-hot-encoding to store the current game state (e.g. cards in each players' hand, cards that have been played or discarded, and what the current player knows from hints) and for the player's action on that turn (e.g. play card 1, discard card 4, etc.). The data is a matrix with dimensions 192731 x 393. The first 373 features represent Game State, the last 20 represent Player Action.

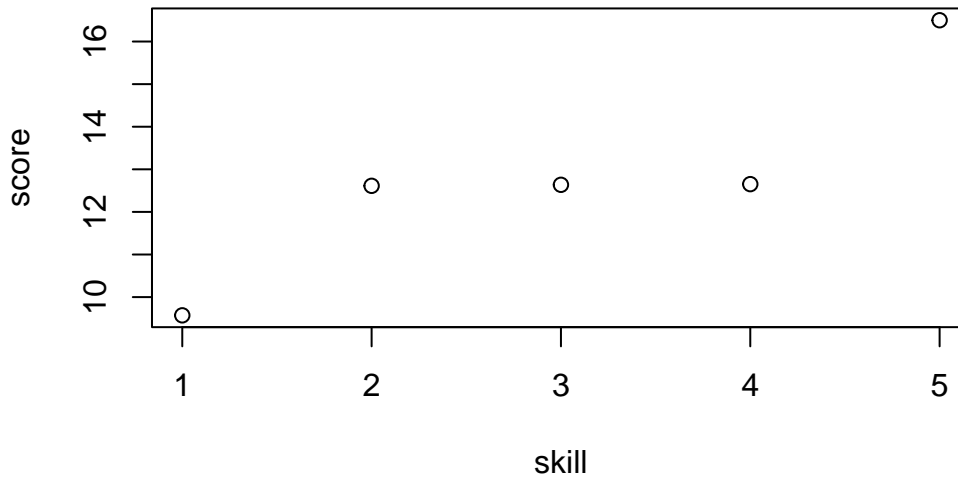
Study Objectives

Our study has three main objectives. First is to reproduce the results from the original paper using the extended dataset. Second is to extend the research to consider player perception of the AI (e.g. skill, likability). Lastly is to train a machine learning model on the game logs to predict the player's next move based on the current game state. This model could then be a proposed addition to the current AI implementation that could be a better collaborator with the human.

Exploratory Data Analysis (1-1.5 pages)

- Summary statistics (as presentable table)
- Data visualization (2-3 key plots)
- Data quality assessment (missing values, outliers, distributions)
- Preliminary insights that inform modeling choices





Track 1 (Paper Reproduction) should include:

- A combination of original EDA and useful pieces of the EDA presented in the paper. If you use/reproduce a plot from the original paper in your report, give credit to the authors in the source.

Methods (1-1.5 pages)

In the original study, participants were told they would be playing with an AI, but not which AI implementation they would be matched with. The original analysis involved an ANOVA comparing score based on AI implementation (Intentional, Outer, or Full). The original study's ANOVA resulted in ($p < 0.0001$) showing statistical significance for a relation between AI implementation and final score. To compare the specific implementations, the original study used a Tukey test to determine that there was a statistically significant difference between the intentional AI and the outer state AI (baseline), as well as between the outer state AI and the full AI, however they did not find statistical significance when comparing the intentional to the full AI.

Beyond reproducing the original ANOVA and Tukey tests using the larger dataset, we extend this research by considering the potential correlation of player perception of AI on the resulting score regardless of AI implementation. To measure this correlation, we use Spearman's Rank considering the 5-point Likert responses to perceived AI Skill and how much the player liked their AI companion as potential variables. This method is used to test how strongly

two sets of ranks are correlated, in this case we hypothesize that both perceived AI skill and likability will be correlated with higher game score.

Lastly, we used Machine Learning to train a model to predict player action based on current game state. Since the game log data was structured and labeled, we used (XGBoost or Neural Net) to train the model.

Results (2-2.5 pages)

```

              Df Sum Sq Mean Sq F value    Pr(>F)
ai              2     589   294.66    11.19 2.42e-05 ***
Residuals    210    5531    26.34
---

```

```

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

Tukey multiple comparisons of means
 95% family-wise confidence level

```

```

Fit: aov(formula = score ~ ai, data = survey_complete)

```

```

$ai
              diff          lwr          upr          p adj
intentional-full  3.1428571  1.064447  5.221268  0.0012846
outer-full        -0.7467532 -2.747293  1.253786  0.6528219
outer-intentional -3.8896104 -5.921657 -1.857564  0.0000309

```

Track 1 (Paper Reproduction) should include:

- Summary of reproduction success/challenges
- What did the original authors do well? What could be improved?
- Summary of results with proposed model following the next section details.

Both tracks should include:

- Model fitting and diagnostics
 - Assumption checking (residual plots, normality tests, etc.)
 - Model comparison (if applicable)
 - Goodness-of-fit measures
- Parameter interpretation with confidence intervals where appropriate

- Key findings presented with visualizations

Discussion & Conclusion (0.75-1 page)

Both tracks should include:

- Answers to research questions with supporting evidence
- Practical implications or insights
- Limitations and assumptions
- Future directions

Track 1 (Reproduction) should include:

- A summary of what could be improved in the original paper and how your new analysis does exactly that.

References

- Eger, Markus, Chris Martens, and Marcela Alfaro Cordoba. 2017. “An Intentional AI for Hanabi.” In *2017 IEEE Conference on Computational Intelligence and Games (CIG)*, 68–75. <https://doi.org/10.1109/CIG.2017.8080417>.
- Hu, Chengpeng, Yunlong Zhao, Ziqi Wang, Haocheng Du, and Jialin Liu. 2024. “Games for Artificial Intelligence Research: A Review and Perspectives.” *IEEE Transactions on Artificial Intelligence* 5 (12): 5949–68. <https://doi.org/10.1109/TAI.2024.3410935>.
- Laird, John E. 2002. “Research in Human-Level AI Using Computer Games.” *Commun. ACM* 45 (1): 32–35. <https://doi.org/10.1145/502269.502290>.
- Yannakakis, Georgios N. 2005. “AI in Computer Games: Generating Interesting Interactive Opponents by the Use of Evolutionary Computation,” December. <https://era.ed.ac.uk/handle/1842/879>.
- Yin, Qi-Yue, Jun Yang, Kai-Qi Huang, Mei-Jing Zhao, Wan-Cheng Ni, Bin Liang, Yan Huang, Shu Wu, and Liang Wang. 2023. “AI in Human-Computer Gaming: Techniques, Challenges and Opportunities.” *Machine Intelligence Research* 20 (3): 299–317. <https://doi.org/10.1007/s11633-022-1384-6>.