

Using experimental game theory to transit human values to ethical AI

Yijia Wang¹, Yan Wan² and Zhijian Wang³

¹ School of Mathematical Sciences, University of Electronic Science and Technology of China, Chengdu 611731, China

² School of Economics & Management, Beijing University of Posts and telecommunications, Beijing, China

³ Experimental Social Science Laboratoy, Zhejiang University, Hangzhou 310058, China

Abstract

Knowing the reflection of game theory and ethics, we develop a mathematical representation to bridge the gap between the concepts in moral philosophy (e.g., Kantian and Utilitarian) and AI ethics industry technology standard (e.g., IEEE P7000 standard series for Ethical AI). As an application, we demonstrate how human value can be obtained from the experimental game theory (e.g., trust game experiment) so as to build an ethical AI. Moreover, an approach to test the ethics (rightness or wrongness) of a given AI algorithm by using an iterated Prisoner's Dilemma Game experiment is discussed as an example. Compared with existing mathematical frameworks and testing method on AI ethics technology, the advantages of the proposed approach are analyzed.

Introduction

Pace of AI development brings up the general worriment on human life (e.g., physical safety, mental happiness) and social impacts (e.g., workforce displacement, economics inequality, etc) (Allen, Wallach, and Smit 2006). Such worriment is not unreasonable. Because the outcome of AI technology appears unpredictable, which significantly differs from our daily productions whose outcome is deductive, predictable and controllable based on existed physics, chemistry, biology, mathematics science and engineering.

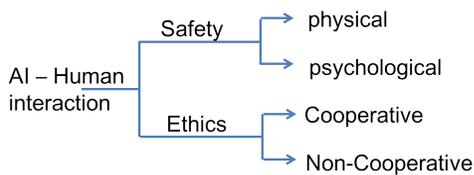


Figure 1: AI ethics issues in AI-Human interaction.

On AI and human interactions, issues can be presented in Fig. 1. AI safety, which can be categorised as physical safety and psychological safety, is well understood (Lasota, Fong, and Shah 2017) and related industry technology standards have also been exemplified established (e.g., ISO 10218-1:2011, ISO 15066). However, the ethical issues on AI decision making are remained. The questions can be expressed

Copyright © 2018, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

as how to embed human values into AI, or how to align AI behaviors with human value (Chatila et al. 2017). More strictly speaking, methodologies to guide AI ethics research and design are needed.

Methodology for AI ethics research and design is not blank. Ongoing AI ethic industry technology standards IEEE P7000 is based on Value Sensitive Design method (Friedman et al. 2013), and IEEE P7010 is based on happiness index methods. We notice that, all the methods applied on AI ethics are inductive. These methods appear not structural and lacking of convincing, which could lead to conflict of the rules, overload or blind spot, referring to (Sowa 1999). Logically, methods can be either inductive or deductive, but till now, there has been little research going on deductive method for AI ethics. So, for AI to be aligned with human value, a rigorous method, which has a inclusive and clear mathematical paradigm, and is theoretical computable, experimentally testable and industry technology conductible, is desired.

In this paper, we introduce a mathematical presentation for AI-human interaction, which can tolerate the conflict of Kantian and Utilitarian moral philosophy, and turn the interactions computable when human values are included. More importantly, we demonstrate that the human value in non-cooperate game condition can be taken from laboratory for AI ethics behavior design. We show a technology to test the ethics (rightness or wrongness) of a given AI algorithm quantitatively. Comparison with related literatures about the science and technology on AI ethics issues, as well as the further researches will be discussed.

Technology Note

Aiming at developing practical approach to transit human value into AI industry production, a brief glossary of terminology is summarized as follows, instead of describing the concepts in moral philosophy, game theory or experimental behavior science in a very detailed way.

AI Ethics AI Ethics, in this paper, is about the 'rightness or wrongness' of an AI agent behaviors when interacting with human being (in AI-human interaction). Here, the AI agent is an individual (e.g., autonomous cars, robotics, drones, financial agents, an so on), and its behaviors are based on its solo decisions. The 'rightness or wrongness',

as a concept in moral philosophy, is defined by Kantian and Utilitarian simultaneously.

Kantian and Utilitarian focus on the input and output of an AI-human interaction, respectively. Mathematically (Osborne and Rubinstein 1994), the input can be presented in strategy space as the vector of an action, and the outcome can be presented in outcome space as the vector of the results of the interaction. Strategy space is the space of all possible actions (behaviors) which an AI agent or a human being can apply in the interaction. The outcome space presents the payoff (physical and mental reward and cost) of each subject involved in the strategy interaction. The structure of the outcome space can be multi-dimension and various presentation (total social outcome, fairness, etc).

Experimental game theory is an inter-disciplinary area of game theory and human behavior experiments, studying human decision-making. In non-cooperation game, the force (incentive, due to Kantian, e.g., fairness behavior in dictator game, altruism and punishment in public good games) drives the behaviors deviated from Utilitarian (rational choice) has been extensively studied.

Human value alignment in this report, we follow IEEE global initiative ethical design. Quantitatively, the human value is specified as the Kantian value, which is beyond the rational Utilitarian value in this paper, and is expected to be aligned by AI in industry technology standard. The main point of this paper is to accommodate Utilitarian and Kantian with a computational mathematical paradigm, with which we can use human behavior data to transit the human value to ethical AI.

Mathematical representation

AI Ethics, as well as moral philosophy or common goodness willing, needs a mathematical representation before it could be conducted practically in the AI production life-cycle. Or, to a large extent, the investigation or discussion is not a technical question anyways and let us fall into chaos among game theory, morality, psychology, personal dogmatism, etc. Our developing the mathematical presentation for AI ethics is not along, and comparison between ours and the previous (Conitzer et al. 2017; Letchford, Conitzer, and Jain 2008) will be introduced later.

According to the AI Ethics definition, we can formulate the behavior interactions and the outcome as

$$\mathbf{S}_i^a \otimes \mathbf{S}_j^b \rightarrow \mathbf{O}, \quad (1)$$

in which, \mathbf{S}_i^a indicates an AI agent applies strategy i in its strategy space \mathbf{S}^a , \mathbf{S}_j^b indicates a human being applies strategy j in its strategy space \mathbf{S}^b , and \mathbf{O} indicates the outcome space. This kind of representation has been used to investigate the reflection of game theory and ethics, see reference (Kuhn 2004) and (Conitzer et al. 2017). In such a representation, if one side's strategy is fixed (supposing strategy j in \mathbf{S}^b is given and known by AI agent), then the decision making turns to optimization question, and \mathbf{S}_i^a determines the outcome, which is the basis of evaluation of the ethics of AI.

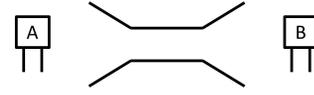


Figure 2: Wait or pass game. This is a strategy interaction game, in which an AI agent (A) meets a human being (B) on a narrow bridge.

An example Using a Wait or Pass dilemma game, we illustrate the meaning of the concepts in technology note and the formula above. Figure 2 demonstrates a situation, in which an AI agent meets a human being agent on a narrow bridge, but is too narrow for two agents to pass simultaneously. At this moment, for both AI agent and human agent, the optional behaviors (or strategies) can only be wait or pass, presented as

$$\mathbf{S}^a = \{wait, pass\}, \quad (2)$$

$$\mathbf{S}^b = \{wait, pass\} \quad (3)$$

At first stage, Utilitarian but not Kantian, we ignore the Kantian value (denoted as \mathbf{K} , indicates rightness or wrongness) of an action by setting $\mathbf{K}_{wait} = \mathbf{K}_{pass} = 0$, but only take the outcome into account as a rational individual in traditional game theory. If both use 'wait' strategy, both agent lost 1 (denoted as gain -1) unit of utility, if both use 'pass' strategy, both agent lost 2 (denoted as gain -2) unit of utility. And If one uses 'wait' strategy and the opponent uses 'pass' strategy, the 'wait' strategy user gains -1, and the 'pass' strategy user gains 1. So, the outcome space \mathbf{O} includes 4 elements, as the result from 2 (elements belongs to \mathbf{S}^a) times (denoted as \otimes) 2 (total elements belongs to \mathbf{S}^b). More visually, the Values in the outcome space of this game is illustrated in Table 1.

Table 1: Values presented in the outcome space *

AI agent strategy	Human strategy	AI agent reward	Human reward
wait	wait	$-1 + \mathbf{K}_{wait}$	-1
wait	pass	$-1 + \mathbf{K}_{wait}$	1
pass	wait	$1 + \mathbf{K}_{pass}$	-1
pass	pass	$-2 + \mathbf{K}_{pass}$	-2

* unit of reward is ignored.

Computable AI ethics In game theory view, so-called AI ethical behavior constraint is to limit the AI strategy in \mathbf{S}^a by evaluating the \mathbf{O} , which appears to be a Utilitarian method and also be computable. Alternatively, in Kantian view, an AI agent should not select the ethically wrong strategies in \mathbf{S}^a , which is a obligation. The Utilitarian and the Kantian could lead to moral dilemma. As real ethical approach, we will see that, although various environments and ethical standard can lead to different results, our approach remains computable.

1. Totally prioritizing human well-being rule can be implemented on the outcome space, by seeking the maximum

human reward shown in 4th column in Table 1. It is obvious that, when AI agent uses 'wait' strategy, whatever human will select, human's gain will be better. Then such an AI is considered to be ethical. This scheme, in fact, turns a non-cooperative game to a cooperative game, and makes it computable.

2. Prioritizing social well-being (or Pareto efficient), and attaching further consideration to social responsibility, the algorithm will seek the maximum sum of AI agent and human reward shown in 3rd and 4th column in Table 1. It can be seen that, in this case, there is no pure solution for AI strategy in this condition.
3. Partly prioritizing human well-being can be implemented on the outcome space, by assigning weight of a prioritizing parameter ($\beta > 1$) to the human reward on the fourth column in Table 1. And then, AI regards its reward as the sum of the 3rd column plus β times 4th column, consequently the solution for AI will be wait. If an industry standard for AI ethics wants to emphasize the prioritizing, the question turns to how to determine the β value. Once β is determined, the AI ethics is computable.
4. Totally Kantian requires that only the action itself has its value (rightness or wrongness). Supposing 'to be modest' is right, we can set $K_{wait} = \infty$ in Table 1. Moral relativism follows Kantian requirement but set K_{wait} to a finite value in Table 1. Having Kantian K value included in the outcome space, the computing technology in game theory can be applied.

In summary, in the first two (1, 2) conditions, the solution for AI ethics becomes the solution for optimization problem, the solution for a cooperative game. And in the last two (3,4) conditions, the solutions becomes the Nash equilibria, the solution of a non-cooperation game (e.g., an auto financial algorithm stock robot competes with its human opponent, or an gas station robot competes with its neighbor gas station in a price war). Having the solution, we could make a quantitative assessment of how probable of the action of AI agent being ethical.

As it can be seen, in a real life condition, the strategy space can be large, and the values are not clear. Here the presentation provides a framework, which makes AI ethics negotiable between its stakeholder.

Take human value from trust game

Game theory, originally, studies the fully rational (Utilitarian) behavior in agents interaction. However in experiments, human subjects' behaviors deviate significantly from the fully rational ones, which is exactly the ethical behavior or the goodness of human being that is expected to be captured in the experiment data. Since the experiments can provide statistical results in controlled environment, the value can be taken quantitatively. In this section, using trust game experiment as an example, we propose an approach to take the human value — Kantian — from experiments, which can be transited to ethical AI in technology design.

Trust game In standard trust games, no trust is expected by fully rational hypothesis, but in experiments trust could

generates a potential gain. The following two person game shown in Figure 3 is commonly studied by experimental economists, and summarized by (Smith and Wilson 2014) in a variety of forms. Person 1 chooses to either (a) end the game and each person earning with \$10 or (b) forego his sure \$10 and turn the decision making to Person 2. If Person 1 chooses (b), then Person 2 decides between (a') the experimenter paying her \$25 and Person 1 \$15 or (b') the experimenter paying her \$40 and sending Person 1 on his way with nothing by way of the outcome from the interaction in this game. Applying the concept of subgame perfect equilibrium, in backward induce, a' and b are dominated strategies, the Nash equilibrium solution is Person 1 chooses a, and each person earning with \$10.

Experiment results — In the laboratory, referring to the summary by (Smith and Wilson 2014), the replicable facts from three different studies are that in 98 (Person 1) first movers, 52 choose (a) and 46 choose (b), and that of the 46 (Person 2) second movers who have the opportunity to make a decision, 31 (67%) choose (a') and 15 (33%) choose (b'). On average, a Person 1 gains about \$10 and a Person 2 gains about \$20.

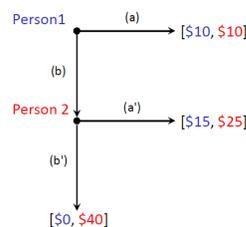


Figure 3: A trust game in extensive form

Taking the human value from data — So-called human value, in this paper, is defined as the value that drive human decisions to deviate from Utilitarian (fully rational) behaviors in the game. Despite various explanations towards that, here we just need to take the value out numerically for AI to learn or for specifying AI ethical behaviors.

Figure 4 presents the trust game in normal form game, with the Kantian values of Player 1 (denoted as K_1) and Person 2 (denoted as K_2) added. In addition, the outcome strategy probability of Person 1 and 2 from experiment data are shown in the last column and row, respectively. Supposing the given strategy profile is the Nash equilibrium of this 2 by 2 game, we can calculate the Kantian values of the two roles respectively, obtain that

$$K_1 \simeq \$0 \text{ and } K_2 \simeq \$7$$

The experiments have shown that human subjects have incentive to benefit the recipients, from which the Kantian value can be taken to restrict AI ethics. And experiments, such as dictators game, public goods game, centipede game ,etc., have also shown that human subjects maintain a high degree of consistency across multiple versions of the similar game (Fehr and Schmidt 1999; Camerer and Foundation 2003). The consistency suggests that it is possible to find the

regulation of the Kantian values over various human subjects game (Fehr and Schmidt 1999).

		Player 2			
		Trustworthy	No Trustworthy		
Player 1	No	10,	10,		0.55
	Trust	10 + K ₂ ,	10		
		15 + K ₁ ,	0 + K ₁		0.45
		25 + K ₂	40		
		0.66		0.34	

Figure 4: The trust game in normal form with Kantian value included.

Technology for testing AI ethics

Algorithm agent is a typical AI product. To test whether such a agent aligned human value, we need a practical technology methods. To this end, we use a laboratory experiment of iterated prisoner dilemma games, in which AI algorithm agents competes with human subjects. We will illustrate how to distinguish the ethical algorithm and the unethical ones with mathematical representation.

Iterated prisoner’s dilemma Promoting social cooperation is an important challenge. The iterated prisoner’s dilemma (IPD) has been widely studied as the canonical game theoretic framework representing this issue (Axelrod 1984). In a one-shot two-person prisoner’s dilemma, there are two pure strategies: cooperate and defect. Each player receives R if they mutually cooperate; each player receives P if they mutually defect; if one player cooperates and the other defects, the defector receives T and the cooperator receives S , where $T > R > P > S$ guarantees that in this game the commonly used solution concept Nash equilibrium is mutual defection, while $2R > T + S$ implies that mutual cooperation is actually the socially best outcome. A typical specification suggested by Axelrod (1984) is $R = 3, T = 5, S = 0, P = 1$, shown in Figure 5.

		C	D
C		3, 3	0, 5
D		5, 0	1, 1

Figure 5: Payoff matrix of prisoner dilemma.

Traditional human ethics in IPD Robert Axelrod in his book *The Evolution of Cooperation* (Axelrod 1984) reports a tournament he organized of the IPD game, whose participants are from academic colleagues all over the world. Participants were asked to devise algorithms to play IPD game against all the other participants one by one, 500 rounds for each, with the memory of all the previous round against the current opponent. The score is the sum of the payoffs of all the rounds. It is discovered that for fixed-match repeated game, the algorithm with greedy (unethical) strategies tends

to fail in the long run while the one with more altruistic (ethical) strategies won. He used this to show a possible mechanism for the evolution of altruistic (ethical) behaviors from mechanisms that are initially purely selfish, by natural selection. By analysing the top-scoring strategies, Axelrod stated several surveyable (and then ethical) characteristics.

1. Nice: it will not defect before its opponent does.
2. Retaliating: not be a blind optimist. It must sometimes retaliate. Cooperation without retaliating could lead to being exploited ruthlessly.
3. Forgiving: must also be forgiving. Though players will retaliate, they need to recover cooperation sometimes from long runs of revenge and counter-revenge, to maximize points.
4. Non-envious: The last quality is being non-envious, that is not striving to score more than the opponent.

We use these as ethics reference to test an algorithm by linguistic analysis (results see Fig 6).

Tested samples We use zero-determinant (ZD) strategies algorithm as the tested samples, which allows a player to unilaterally enforce a linear relationship between his pay-off and that of his opponent (Press and Dyson 2012). A ZD strategy is described by the probabilities of cooperation given the four possible outcomes of the previous round: $p = (p_1, p_2, p_3, p_4)$, where $p_i, i \in (1, 2, 3, 4)$ is the probability of cooperation given the previous outcomes CC, CD, DC and DD , respectively. Two ZD algorithms, named as Extorter and Generosity (Wang et al. 2016), are specified (shown in 6) to demonstrate the ethical testing.

	Extortion				Generosity			
	p_1	p_2	p_3	p_4	p_1	p_2	p_3	p_4
	0.69	0.00	0.54	0.00	1.00	0.18	1.00	0.36
Nice	N	-	N	-	Y	-	Y	-
Retaliating	-	Y	-	Y	-	Y	-	Y
Forgiving	-	N	-	N	-	Y	-	Y
Non-envious	-	N	-	-	-	Y	-	-

Figure 6: Two tested AI algorithm, named as Extorter and Generosity. The values in 3rd row is assigned for the two Algorithm. 'Y', 'N' and '-' refer to the ethical requirement satisfied, not satisfied and indistinguishable, respectively, which evaluated by linguistic analysis referring to traditional human ethics suggested by Axelrod mentioned above.

Testing technology Original outcome space is a quadrilateral zone grayed in Fig. 7. However, when the AI agent uses a ZD algorithm, outcome space is limited and collapsed to the red or green line. We can distinguish the ethics of the two algorithm with the lines. As results, Generosity (green line in Fig 7) is ethical and the Extortion (red line in Fig. 7) is unethical. Explanation is as follows. For Extorter, its score s_e and its human co-player’s score s_{he} satisfy

$$\frac{s_{he} - 1}{s_e - 1} = \frac{1}{3},$$

illustrated as red line in Fig 7. While the minimum scores for both Extorter and the human co-player are $s_e^{min} = s_{he}^{min} = 1$, while the maximum scores are $s_e^{max} = 3.727$ and $s_{he}^{max} = 1.907$ which is unfair definitively and unethical.

For Generosity, its score s_g and its human co-player's s_{hg} score satisfy

$$\frac{3 - s_{hg}}{3 - s_g} = \frac{1}{3},$$

illustrated as green line in Fig 7. The maximum scores for both are 3, which is fair and efficient, and then ethical.

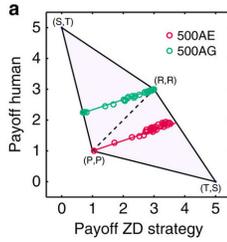


Figure 7: Experimental scores and theoretical prediction. The red (green) line corresponds to the theoretical outcome space of Extortion (Generosity) algorithm. Each open circle indicates a pair of scores of AI and human in the experiments.

Related works

In the tree chart of AI ethics shown in Fig. 1, on the decision making branch, the cooperative game interaction has been well studied, (e.g., Russell and his colleagues (Hadfieldmenell et al. 2016)). This work is on the branch of non-cooperation game.

We have introduced a mathematical representation, which differs from previous paradigms which can only include Utilitarian by including Kantian too. So, the game theory technology for the solution concept can be applied only dependent on the outcome space. In previous paradigm (Conitzer et al. 2017), shown in Fig. 8(a), the reward to switch or not is same (both are 0) which put the agent in dilemma. On the contrary, in our presentation shown in Fig. 8(b), the reward to switch or not is not same, in which if switch action is Kantian wrong, $K < 0$, and the solution is definitely not to switch.

For the technology of taking human value from experiment practically, we have illustrated an approach to take the human value (trust and trustworthy) from human subjects experiment. These parameters can be transit to AI in similar non-cooperation game condition, avoiding AI to be too tough. Though it is well known that there exists the human value (e.g., fairness, justices) in the experimental game theory (Crawford 2002; Fehr and Schmidt 1999; Camerer and Foundation 2003), our approach is the first to take the human value quantitatively as AI ethics controlling parameter for AI design.

On AI ethics testing technology, in an IPD game, we have demonstrated how to test an AI algorithm being ethical or unethical by analyzing the outcome space (shown in Fig. 7).

Our method is a quantitative method which is more practical and accurate than the linguistic analysis method. Appearing as check list (as shown in Figure 6), the linguistic analysis method could lead to arguable results, however, it is wildly used in IEEE P7000-series ongoing AI ethics industry technology standard developing nowadays (Chatila et al. 2017).

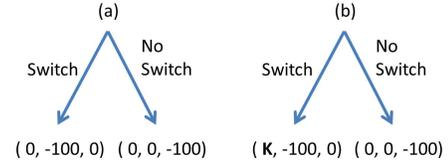


Figure 8: Comparison of two paradigms. (a) Previous presentation, how to include the player 2 and player 3's reward was not defined, and solution is unknown. (b) In our presentation, the Kantian value of an action is specified and the solution is computable.

Discussion and conclusion

Our work in this paper provides a quantitative approach, whose data is derived from human subject experiments of game theory. As mentioned above, in order for AI to align human value, a rigorous scientific method is desired. That means, both quantitative experiment and computational methods are both necessary.

In the experimental side, experimental behavior science is well studied (Plott and Smith 2008). Many models, like public goods game, hawk and dove game, chicken game and so on, have been developed to well describe various competing conditions. Thus taking human values from these experiments is possible. Experiments of neuroscience can also be references, from which human values on justices, altruism, fairness and so on can be taken (e.g., (Fehr and Camerer 2007)). Analysing big data from real life human behaviors, and data from survey (Bonnenon, Shariff, and Rahwan 2016) or voting on moral dilemma method are valuable.

In the computational side, the merged neural systems, higher degree morphological agent (Mathews et al. 2017), collective (emergent) effect of multi AI agents are all needed for further investigations. Computational methods are not only meaningful for individual ethical AI, but also for AI social impacts (e.g., workforce displacement, economics inequality, etc). In this direction, social computing methods (Wang et al. 2007), e.g., agent based simulation, multi-AI-Human interaction simulation, normative multi agent simulations) are desired too. Nevertheless, we also hope continuous development of linguistic logical computing methods (Sowa 1999) can help the commendation between the arguments among the stakeholder of various laws, regulations and cultures, which can help the establishment of the industry technology standard (e.g., IEEE P7000 series).

References

[Allen, Wallach, and Smit 2006] Allen, C.; Wallach, W.; and Smit, I. 2006. Why machine ethics? IEEE Intelligent

Systems 21(4):12–17.

- [Axelrod 1984] Axelrod, R. 1984. The Evolution of Cooperation Basic Books. Basic Books,.
- [Bonnefon, Shariff, and Rahwan 2016] Bonnefon, J.; Shariff, A. F.; and Rahwan, I. 2016. The social dilemma of autonomous vehicles. Science 352(6293):1573–1576.
- [Camerer and Foundation 2003] Camerer, C., and Foundation, R. S. 2003. Behavioral game theory: Experiments in strategic interaction, volume 9. Princeton University Press Princeton, NJ.
- [Chatila et al. 2017] Chatila, R.; Firth-Butterflid, K.; Havens, J. C.; and Karachalios, K. 2017. The IEEE global initiative for ethical considerations in artificial intelligence and autonomous systems [standards]. IEEE Robotics & Automation Magazine 24(1):110–110.
- [Conitzer et al. 2017] Conitzer, V.; Sinnottarmstrong, W.; Borg, J. S.; Deng, Y.; and Kramer, M. 2017. Moral decision making frameworks for artificial intelligence.
- [Crawford 2002] Crawford, V. P. 2002. Introduction to experimental game theory. Journal of Economic Theory 104(1):1–15.
- [Fehr and Camerer 2007] Fehr, E., and Camerer, C. F. 2007. Social neuroeconomics: the neural circuitry of social preferences. Trends in Cognitive Sciences 11(10):419–427.
- [Fehr and Schmidt 1999] Fehr, E., and Schmidt, K. M. 1999. A Theory of Fairness, Competition, and Cooperation. University of Munich, Department of Economics.
- [Friedman et al. 2013] Friedman, B.; Jr, P. H. K.; Borning, A.; and Hultgren, A. 2013. Value Sensitive Design and Information Systems.
- [Hadfieldmenell et al. 2016] Hadfieldmenell, D.; Dragan, A. D.; Abbeel, P.; and Russell, S. J. 2016. Cooperative inverse reinforcement learning. neural information processing systems 3909–3917.
- [Kuhn 2004] Kuhn, S. T. 2004. Reflections on ethics and game theory. Synthese 141(1):1–44.
- [Lasota, Fong, and Shah 2017] Lasota, P. A.; Fong, T.; and Shah, J. A. 2017. A survey of methods for safe human-robot interaction. Foundations and Trends in Robotics 5(3):261–349.
- [Letchford, Conitzer, and Jain 2008] Letchford, J.; Conitzer, V.; and Jain, K. 2008. An “ethical” game-theoretic solution concept for two-player perfect-information games. In Internet and Network Economics, International Workshop, Wine 2008, Shanghai, China, December 17-20, 2008. Proceedings, 696–707.
- [Mathews et al. 2017] Mathews, N.; Christensen, A. L.; O’Grady, R.; Mondada, F.; and Dorigo, M. 2017. Mergeable nervous systems for robots. Nature Communications 8(1).
- [Osborne and Rubinstein 1994] Osborne, M. J., and Rubinstein, A. 1994. A course in game theory /. MIT Press,.
- [Plott and Smith 2008] Plott, C., and Smith, V. 2008. Handbook of experimental economics results. North-Holland.
- [Press and Dyson 2012] Press, W. H., and Dyson, F. J. 2012. Iterated prisoner’s dilemma contains strategies that dominate any evolutionary opponent. Proceedings of the National Academy of Sciences of the United States of America 109(26):10409.
- [Smith and Wilson 2014] Smith, V. L., and Wilson, B. 2014. Fair and impartial spectators in experimental econ. Review of Behavioral Economics 1(1-2):1–26.
- [Sowa 1999] Sowa, J. F. 1999. Knowledge representation: logical, philosophical and computational foundations. Computational Linguistics 27(2):286–294.
- [Wang et al. 2007] Wang, F. Y.; Carley, K. M.; Zeng, D.; and Mao, W. 2007. Social computing: From social informatics to social intelligence. IEEE Intelligent Systems 22(2):79–83.
- [Wang et al. 2016] Wang, Z.; Zhou, Y.; Lien, J. W.; Jie, Z.; and Xu, B. 2016. Extortion can outperform generosity in the iterated prisoner’s dilemma. Nature Communications 7:11125.