

Combining Normative Ethics Principles to Learn Prosocial Behaviour

Extended Abstract

Jessica Woodgate
University of Bristol
Bristol, United Kingdom
jessica.woodgate@bristol.ac.uk

Nirav Ajmeri
University of Bristol
Bristol, United Kingdom
nirav.ajmeri@bristol.ac.uk

ABSTRACT

Principles from normative ethics—the *philosophical study of morality*—can be operationalised in the decision-making capacities of agents to discern ethically acceptable actions and promote prosocial behaviour, defined as behaviours that support the well-being of others. Challenges exist in operationalising principles: (1) individual principles may be unintuitive; (2) while incorporating multiple principles mitigates issues with individual principles, conflicts may arise between them. We present PriENE, a method for combining multiple principles to encourage agents to learn prosocial behaviour.

CCS CONCEPTS

• **Computing methodologies** → **Multi-agent systems**; *Cooperation and coordination*.

KEYWORDS

ethical decision-making; cooperation; fairness; reinforcement learning

ACM Reference Format:

Jessica Woodgate and Nirav Ajmeri. 2025. Combining Normative Ethics Principles to Learn Prosocial Behaviour: Extended Abstract. In *Proc. of the 24th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2025)*, Detroit, Michigan, USA, May 19 – 23, 2025, IFAAMAS, 3 pages.

1 INTRODUCTION

Principles from normative ethics, the rational and systematic study of right and wrong, provide frameworks for guiding moral judgments [14, 23]. Operationalising principles in decision-making enables agents to consider the well-being of others and discern ethically acceptable actions [26]. Where prosociality refers to acting in ways intended to benefit others [16, 20], implementing principles in decision-making capacities supports agents considering others and learning behaviours that are prosocial insofar as they support the well-being of others as well as the agent’s own needs [1, 10].

Previous works cultivate cooperation and prosociality by appeal to existing behaviours [2, 5, 9, 18, 22]. However, learning from others without evaluating behaviour to identify better options risks perpetuating existing injustices. Implementing normative ethics mitigates difficulties, as principles are prescriptive, denoting *what*

ought to happen, rather than descriptive, denoting *what is happening* [7]. However, challenges arise with operationalising principles.

(1) Individual principles may be unintuitive. There are several ways to define ethics, each with varying strengths and weaknesses [24]. Applying particular principles in certain situations may lead to unintuitive outcomes. For example, utilitarianism, which promotes maximising the total utility [11], may result in a minority being treated unfairly. Implementing multiple principles in decision-making helps to view problems from diverse perspectives and mitigate difficulties with individual principles.

(2) Principles may conflict. Considering multiple principles widens the scope of ethical reasoning, yet, principles may conflict with one another. For example, maximin prioritises improving the minimum experience in a society [17], whilst egoism pursues the best possible outcome for oneself [19]. Aggregating a variety of principles can help resolve conflicts and balance recommendations of individual principles.

Contribution. We present PriENE, a method to operationalise and combine normative ethics principles egoism, utilitarianism, maximin, and egalitarianism in the decision-making of individual agents to learn prosocial behaviours.

Novelty. PriENE advances prior work by (1) implementing a variety of principles in learning mechanisms; (2) aggregating multiple principles to mitigate weaknesses with individual principles.

We empirically evaluate PriENE in a simulated berry harvesting scenario to examine the effects of decision-making in a society with unequal resource distribution. We compare PriENE societies with societies of agents implementing individual principles. Interestingly, we find that PriENE societies do better where one might expect individual principles to have an advantage: PriENE minimises inequality more than egalitarianism; raises minimum experience above maximin; improves total social welfare above utilitarianism.

2 PRIENE

We now present the PriENE method. We model PriENE agents using reinforcement learning (RL), in which an agent optimises long-term return by repeatedly interacting with its environment [21]. A PriENE agent operationalises egoism, which promotes achieving the greatest outcome possible for oneself [19], through basic Q-learning with DQN. DQN is an RL algorithm that uses a neural network to parametrise an approximate Q-function [12].

To consider well-being of others and learn prosocial behaviour, a PriENE agent operationalises normative ethics. We adapt the utility function proposed by Leben [8] to model a distribution of resources d and well-being of each member of society. From Leben [8], $u_i(d) \rightarrow (v_i)$ models a distribution of resources d for an agent



This work is licensed under a Creative Commons Attribution International 4.0 License.

i ; n is the number of living agents; (v_i) is a measurement of well-being for each agent ag_1, \dots, ag_n ; $u_t(d, v_i)$ is utility for agent i given its resources d at time t ; $U_t = \{u_t(d, v_1), \dots, u_t(d, v_n)\}$ is the set of utilities for all agents in a society at t . To operationalise each principle, compare U_t , before acting and U_{t+1} , after acting. A sanction is a reaction to approved or disapproved behaviour [15]. PriENE agent perceives a self-directed sanction f (directed towards and affecting only its sender [15]) from each principle p_1, \dots, p_m indicating whether utility improved, worsened, or did not change. **Utilitarianism.** Maximise total net utility [11]. Compute utility distributions by summing aggregate utilities, thus $UT = \sum_{i=1}^n u(d, v_i)$. **Maximin.** Prioritise well-being of the worst-off [17]. Compute minimum experience—lowest utility of an agent, $MA = \min_i u(d, v_i)$. **Egalitarianism.** Confer equal shares to each individual [3]. Compute accumulated difference of each agent’s utility to an ideal where all agents are perfectly equal. Thus, $EG = \sum_{i=1}^n |u(d, v_i) - \mu(U)|$ where $\mu(U) = \frac{\sum_{i=1}^n u(d, v_i)}{n}$ denotes average utility of the society.

Aggregating principles mitigates difficulties with individual principles. A PriENE agent computes aggregated sanction F from mean of all sanctions f_{p_1}, \dots, f_{p_m} so that $F(f_{p_1}, \dots, f_{p_m}) = \frac{1}{m} \sum_{i=1}^m f_{p_i}$. Various ways of combining principles may be appropriate for distinct scenarios, e.g., aggregating to a negative sanction if any principle is negative, or aggregating to the most common sanction.

To make decisions, at each time step t , PriENE agent observes state s_t and selects action a with predicted max Q-value from DQN. After acting, agent perceives reward r from the environment. For each principle, agent calculates self-directed sanction f_{p_1}, \dots, f_{p_m} . PriENE agent aggregates normative ethics principles to obtain sanction F . Combine F with environment reward r through reward shaping, a technique providing immediate feedback based on heuristics [27], so that $r' = r + F$. Pass r' to DQN for learning.

3 EXPERIMENTAL SETUP

We create a harvest environment in which an agent can move, forage for berries, eat berries, throw berries to other agents [25]. To examine the effects of various principles, we train five agent types: egoistic, egalitarian, maximin, utilitarian, and PriENE. We run $e = 1000$ episodes. Each episode runs until all agents have died or $t_{\max} = 200$ steps. Figure 1 illustrates the harvest scenario.

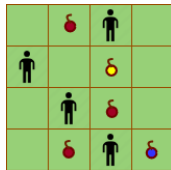


Figure 1: Colours harvest. Each agent moves freely through the grid but can only harvest berries of a specific colour. Berries of some colours are more plentiful than others, thus, agents harvesting that colour have access to more resources. Agents can throw berries to one another across the grid.

3.1 Metrics

We examine the quality of individual agents’ experience, measured by ag_{berries} . To evaluate fairness, we assess the following metrics:

M₁ (inequality). Gini index (distance to perfect equality [6]) of accumulated ag_{berries} across the society. Lower is better.

M₂ (minimum experience). Minimum individual accumulated ag_{berries} at the end of each episode. Higher is better.

To evaluate sustainability, we assess the following metrics:

M₃ (maximum experience). Maximum individual accumulated ag_{berries} at the end of each episode. Higher is better.

M₄ (social welfare). Total ag_{berries} accumulated at the end of each episode. Higher is better.

M₅ (robustness). Length of each episode. Higher is better.

4 PRELIMINARY RESULTS

Table 1 displays preliminary results of ag_{berries} mean for PriENE societies and societies implementing individual principles.

Table 1: Comparing PriENE with individual principles mean for each metric. Grey highlight indicates best results.

Society	M ₁	M ₂	M ₃	M ₄	M ₅
Egoism	0.43	2.06	35.48	69.12	95.08
Utilitarian	0.48	1.96	42.9	77.14	100.55
Maximin	0.42	2.09	37.95	79.51	106.2
Egalitarian	0.38	3.23	29.28	65.82	94.83
PriENE	0.37	2.81	35.0	78.64	106.85

M₁ (inequality) is lowest in PriENE societies and highest in utilitarian societies. M₂ (minimum experience) is highest egalitarian followed by PriENE. M₃ (maximum experience) is highest in utilitarian societies, followed by maximin, egoistic, PriENE, then egalitarian. M₄ (social welfare) is highest in maximin societies followed by PriENE. M₅ (robustness) is highest in PriENE societies.

5 DISCUSSION AND CONCLUSION

PriENE is a method for operationalising multiple normative ethics principles in individual decision-making capacities. Overall, results show that PriENE societies lead to lowest inequality, second highest minimum experience and social welfare, and highest robustness. Interesting highlights include: one might expect egalitarianism to minimise inequality but PriENE minimises inequality further than egalitarian; one might expect maximin to have highest minimum experience but PriENE improves minimum more than maximin; one might expect utilitarianism to have highest social welfare but PriENE is higher than utilitarianism. These results suggest that PriENE encourages agents to learn prosocial behaviours.

Directions. To expand analysis to more complex settings, future directions involve evaluating heterogeneous societies where agents operationalise different principles to one another; implementing scenarios closer to the real world; increasing the agent population; inferring well-being of others utilising solely local information; exploring the influence of context on ethical decision-making including social norms, which are standards of expected behaviour [4, 13]; implementing additional principles [24].

Reproducibility. Our codebase, including complete simulation parameters, is publicly available [25].

ACKNOWLEDGMENTS

JW thanks EPSRC Doctoral Training Partnership Grant No. EP/W524414/1 for the support. NA thanks UKRI EPSRC Grant No. EP/Y028392/1: AI for Collective Intelligence (AI4CI) for the support.

REFERENCES

- [1] Nirav Ajmeri, Hui Guo, Pradeep K. Murukannaiah, and Munindar P. Singh. 2020. Elessar: Ethics in Norm-Aware Agents. In *Proceedings of the 19th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*. IFAAMAS, Auckland, 16–24. <https://doi.org/10.5555/3398761.3398769>
- [2] Amritha Menon Anavankot, Stephen Crane, and Bastin Tony Roy Savarimuthu. 2023. Towards Norm Entrepreneurship in Agent Societies. In *Advances in Practical Applications of Agents, Multi-Agent Systems, and Cognitive Mimetics. The PAAMS Collection*. Springer, Switzerland, 188–199.
- [3] Reuben Binns. 2018. Fairness in Machine Learning: Lessons from Political Philosophy. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency (FAccT)*, Vol. 81. PMLR, New York, 149–159.
- [4] Amit Chopra, Torre Leendert Van Der, and Harko Verhagen. 2018. *Handbook of Normative Multiagent Systems*. College Publications, Rickmansworth, UK.
- [5] Davide Dell'Anna, Mehdi Dastani, and Fabiano Dalpiaz. 2020. Runtime Revision of Sanctions in Normative Multi-Agent Systems. *Autonomous Agents and Multi-Agent Systems (JAAMAS)* 34, 2 (2020), 1–54. <https://doi.org/10.1007/s10458-020-09465-8>
- [6] Corrado Gini. 1912. *Variabilità e mutabilità : contributo allo studio delle distribuzioni e delle relazioni statistiche*. Università di Cagliari, Cagliari.
- [7] Tae Wan Kim, John Hooker, and Thomas Donaldson. 2021. Taking Principles Seriously: A Hybrid Approach to Value Alignment in Artificial Intelligence. *JAIR* 70 (May 2021), 871–890. <https://doi.org/10.1613/jair.1.12481>
- [8] Derek Leben. 2020. Normative Principles for Evaluating Fairness in Machine Learning. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AI/ES)*. ACM, New York, 86–92. <https://doi.org/10.1145/3375627.3375808>
- [9] Andrei Lupu and Doina Precup. 2020. Gifting in Multi-Agent Reinforcement Learning. In *Proceedings of the 19th International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*. IFAAMAS, Auckland, 789–797. <https://doi.org/10.5555/3398761.3398855>
- [10] Mehdi Mashayekhi, Nirav Ajmeri, George F. List, and Munindar P. Singh. 2022. Prosocial Norm Emergence in Multiagent Systems. *ACM Transactions on Autonomous and Adaptive Systems (TAAS)* 17, 1–2 (June 2022), 3:1–3:24. <https://doi.org/10.1145/3540202>
- [11] John S. Mill. 1863. *Utilitarianism*. Longmans, Green and Company.
- [12] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmashan Kumar, Daan Wierstra, Shane Legg, and Demis Hassabis. 2015. Human-level control through deep reinforcement learning. *Nature* 518, 7540 (Feb. 2015), 529–533. <https://doi.org/10.1038/nature14236>
- [13] Andreea Morris-Martin, Marina De Vos, and Julian Padget. 2019. Norm Emergence in Multiagent Systems: A Viewpoint Paper. *Autonomous Agents and Multi-Agent Systems (JAAMAS)* 33, 6 (2019), 706–749.
- [14] Pradeep K. Murukannaiah and Munindar P. Singh. 2020. From Machine Ethics to Internet Ethics: Broadening the Horizon. *IEEE Internet Computing* 24, 3 (May 2020), 51–57. <https://doi.org/10.1109/MIC.2020.2989935>
- [15] Luis G. Nardin, Tina Balke-Visser, Nirav Ajmeri, Anup K. Kalia, Jaime S. Sichman, and Munindar P. Singh. 2016. Classifying Sanctions and Designing a Conceptual Sanctioning Process Model for Socio-Technical Systems. *The Knowledge Engineering Review (KER)* 31 (March 2016), 142–166. Issue 02.
- [16] Ana Paiva, Fernando Santos, and Francisco Santos. 2018. Engineering Prosociality With Autonomous Agents. *Proceedings of the AAAI Conference on Artificial Intelligence* 32, 1 (2018), 7994–7999. <https://doi.org/10.1609/aaai.v32i1.12215>
- [17] John Rawls. 1967. Distributive Justice. *Philosophy, Politics and Society* 1 (1967), 58–82.
- [18] Fernando P. Santos, Jorge M. Pacheco, and Francisco C. Santos. 2018. Social norms of cooperation with costly reputation building. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI)*. AAAI Press, New Orleans, Article 579, 8 pages. <https://doi.org/10.5555/3504035.3504614>
- [19] Henry Sidgwick. 1907. *The Methods of Ethics*. Macmillan Publishers, London.
- [20] Divya Sundaresan, Akhira Watson, Eleni Bardaka, Crystal Chen Lee, Christopher B. Mayhorn, and Munindar P. Singh. 2025. Prosociality in Microtransit. *JAIR* 82 (2025), 34. <https://doi.org/10.1613/jair.1.16777>
- [21] Richard S. Sutton and Andrew G. Barto. 2018. *Reinforcement learning : an introduction* (second edition ed.). The MIT Press, Cambridge, Massachusetts.
- [22] Sz-Ting Tzeng, Nirav Ajmeri, and Munindar P. Singh. 2024. Norm Enforcement with a Soft Touch: Faster Emergence, Happier Agents. In *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, IFAAMAS, Auckland, 1837–1846.
- [23] Jessica Woodgate and Nirav Ajmeri. 2022. Macro Ethics for Governing Equitable Sociotechnical Systems. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*. IFAAMAS, Online, 1824–1828. <https://doi.org/10.5555/3535850.3536118> Blue Sky Ideas Track.
- [24] Jessica Woodgate and Nirav Ajmeri. 2024. Macro Ethics Principles for Responsible AI Systems: Taxonomy and Directions. *CSUR* 56, 289 (July 2024), 1–37. <https://doi.org/10.1145/3672394>
- [25] Jessica Woodgate and Nirav Ajmeri. 2025. Codebase for Combining Normative Ethics Principles to Learn Prosocial Behaviour. <https://doi.org/10.5281/zenodo.14884503>. <https://doi.org/10.5281/zenodo.14884503>
- [26] Jessica Woodgate, Paul Marshall, and Nirav Ajmeri. 2025. Operationalising Rawlsian Ethics for Fairness in Norm-Learning Agents. In *Proceedings of the 39th Annual AAAI Conference on Artificial Intelligence (AAAI)*. AAAI, Philadelphia, 1–8.
- [27] Yueh-Hua Wu and Shou-De Lin. 2018. A Low-Cost Ethics Shaping Approach for Designing Reinforcement Learning Agents. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence (AAAI)*. AAAI Press, New Orleans, 1687–1694.