

# Occam's razor is insufficient to infer the preferences of irrational agents



Stuart Armstrong and Sören Mindermann (equal contributions)

Future of Humanity Institute, Oxford University



## Summary

Algorithms that infer human preferences from the human policy usually assume that the policy  $\pi$  is (noisily) rational. Since humans are *systematically* irrational, this can lead to catastrophically wrong inferences. Can we ever hope to learn human preferences on tasks that aren't trivial for humans? The only available answer is to learn the human rationality model (planner)  $\hat{p}$  and jointly infer their reward function  $\hat{R}$ .

First, a No-Free-Lunch (NFL) result shows that this is impossible without assumptions about human rationality or reward. NFL results are usually solved by Occam's razor priors. Second, this strategy unexpectedly fails because some nonsensical interpretation of human rationality and reward are very simple.

We will need to investigate which additional assumptions humans, who are also subject to this problem, are making.

## Notation

Let  $\Pi$  be the set of (deterministic) policies. Define  $\mathcal{R} = \{R : \mathcal{S} \times \mathcal{A} \rightarrow [-1, 1]\}$ , the reward functions,  $\mathcal{P} = \{p : \mathcal{R} \rightarrow \Pi\}$  the planners.

## Definition

The pair  $(p, R)$  is *compatible* with  $\pi$  iff  $p(R) = \pi$ .

If  $p(R) = \pi$  and the human follows  $\pi$ ,  $(p, R)$  predicts the human's actions perfectly. One can't update away from  $(p, R)$  by observation, so we need priors.

## Theorem

For all  $R$ , and all  $p$  with  $\pi \in \text{Im}(p)$ , there exists  $p'$  and  $R'$  with  $p'(R) = p(R) = \pi$ . Regret from choosing the wrong  $(p, R)$  can be very large [1].

## Degenerate pairs and rewards

We construct 3 pairs  $(p, R)$  that are clearly undesirable, but valid decompositions of the human policy, i.e.  $p(R) = \pi$ . Then we show these are programs with almost minimal description length.

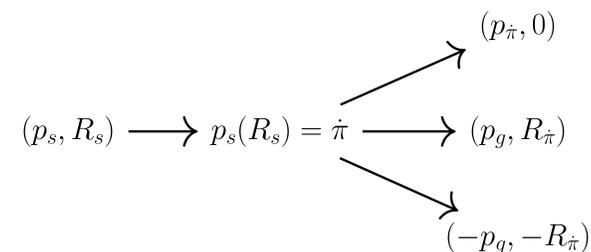
Define  $-p$ :  $-p(R) = p(-R)$ . Note that  $(p, R) \rightarrow (-p, -R)$  is a simple map and  $-p(-R) = p(R)$ . **Thus if  $(\hat{p}, \hat{R})$  is compatible, so is  $(-\hat{p}, -\hat{R})$ .**

For  $\pi$ , define  $p_\pi$ ,  $p_g$  (the greedy planner), and  $R_\pi$ :

$$\begin{aligned} \forall R, s : p_\pi(R)(s) &= \pi(s) \\ p_g(R)(s) &= \text{argmax}_a R(s, a) \\ \forall a, s : R_\pi(s, a) &= \mathbf{1}(\pi(s) = a) \end{aligned}$$

Note that  $(p_\pi, 0)$ ,  $(p_g, R_\pi)$ , and  $(-p_g, -R_\pi)$  are compatible with  $\pi$ . Let  $(p_s, R_s)$  be the simplest pair compatible with  $\pi$ .

Then the following maps are simple programs in any reasonable language  $L$ :



## Complexity of human rationality

Assume  $(\hat{p}, \hat{R})$  fits with human judgement about human rationality and reward. The following imply it is of higher complexity than  $\pi$ :

- 1 Any reasonable  $(\hat{p}, \hat{R})$  is of high complexity, higher than it may intuitively seem to us [2, 3].

- 2 Even given  $\pi$ , any reasonable  $(\hat{p}, \hat{R})$  involves a high number of contingent choices:
  - i Distinct cultures have distinct rationality judgements [4].
  - ii Different individuals have distinct rationality judgements [5].
  - iii The same individual has different judgements at different times [5].
  - iv The judgement of an individual can be changed or manipulated [6].
- 3 Hence any given  $(\hat{p}, \hat{R})$  has extra information (and thus extra complexity), even given  $\pi$ . Points iii and iv prevent 'simply' using the human's own judgement of their rationality.
- 4 Past and present failures to find a simple  $(\hat{p}, \hat{R})$  derived from  $\pi$  are evidence that this is tricky.

## Conclusion

The simplest high-likelihood hypotheses about reward and planner are degenerate. How do humans themselves deal with this problem? They necessarily use some shared priors, that have never been informed by observation. Too strong assumptions can lead to catastrophic false inferences so we need to find minimal sufficient ones.

One potential approach is to look at the 'variables' and models within the human thought process, see Appendix C of the paper.

## References

- [1] T. Everitt, V. Krakovna, L. Orseau, M. Hutter, and S. Legg, "Reinforcement learning with a corrupted reward channel," *arXiv preprint arXiv:1705.08417*, 2017.
- [2] P. W. Glimcher, C. F. Camerer, E. Fehr, and R. A. Poldrack, "Neuroeconomics: Decision making and the brain," 2009.
- [3] L. Muehlhauser and L. Helm, "The singularity and machine ethics," in *Singularity Hypotheses*, pp. 101–126, Springer, 2012.
- [4] J. G. Miller, "Culture and the development of everyday social explanation," *Journal of personality and social psychology*, vol. 46, no. 5, p. 961, 1984.
- [5] P. Slovic and A. Tversky, "Who accepts savage's axiom?," *Behavioral science*, vol. 19, no. 6, pp. 368–373, 1974.
- [6] H. Schuman and J. Ludwig, "The norm of even-handedness in surveys as in life," *American Sociological Review*, pp. 112–120, 1983.

## Simplicity of degenerate planner-rewards and negative rewards

For reasonable languages  $L$ , there are degenerate rationality-reward pairs  $(p, R)$  of minimal Kolmogorov complexity  $K_L$  among those compatible with the human policy  $\pi$ .

$$K_L \left( \begin{array}{c} (p_\pi, 0) \\ (p_g, R_\pi) \\ (-p_g, -R_\pi) \end{array} \right) \approx \min_{(p, R): p(R)=\pi} K_L(p, R) \approx K_L(\pi).$$

Any reasonable human planner-reward pair  $(\hat{p}, \hat{R})$  is of non-negligible higher complexity than these.

$$K_L(\hat{p}, \hat{R}) \gg K_L \left( \begin{array}{c} (p_\pi, 0) \\ (p_g, R_\pi) \\ (-p_g, -R_\pi) \end{array} \right).$$

Also, there is a planner  $p' = -\hat{p}$  such that the complexity of the opposite reward  $-\hat{R}$  is comparable.

$$K_L(\hat{p}, \hat{R}) \approx K_L(-\hat{p}, -\hat{R}).$$

## Contact Information

- Web: [www.fhi.ox.ac.uk/](http://www.fhi.ox.ac.uk/)
- Email: [soeren.mindermann@gmail.com](mailto:soeren.mindermann@gmail.com), [stuart.armstrong@philosophy.ox.ac.uk](mailto:stuart.armstrong@philosophy.ox.ac.uk)