

# Predicting and Preferring

Nathaniel Sharadin ([natesharadin@gmail.com](mailto:natesharadin@gmail.com))  
University of Hong Kong, Department of Philosophy

Forthcoming at *Inquiry*. Please cite published version.

*This work was supported by a grant from the Hong Kong Research Grants Council  
(RGC Grant #17602622)*

## Abstract

The use of machine learning, or “artificial intelligence” (AI) in medicine is widespread and growing. In this paper, I focus on a specific proposed clinical application of AI: using models to predict incapacitated patients’ treatment preferences. Drawing on results from machine learning, I argue this proposal faces a special moral problem. Machine learning researchers owe us assurance on this front before experimental research can proceed. In my conclusion I connect this concern to broader issues in AI safety.

## 1 Introduction: Patient Preference Predictors

There’s a convergence of a long-standing problem in clinical medicine and newly maturing capabilities of predictive models trained using machine learning (ML) -- what are sometimes called “artificial intelligences” (AIs).<sup>1</sup> I’ll describe the problem first, then I’ll describe why it’s natural to suggest using ML, or AI, to solve it.

The problem is simple: care ought to reflect patients’ preferences, but incapacitated patients cannot indicate their preferences. Clinicians can hope to use *indirect* indicators of patients’ preferences, e.g., advance directives and surrogates.<sup>2</sup> But these indirect strategies face serious challenges: most patients do not in fact have an advance directive, and surrogates are systematically epistemically unreliable.<sup>3</sup>

---

<sup>1</sup> In what follows, I mostly refer to these systems as ML systems, rather than as “AI”, in order to avoid unfortunate and controversial implications about machine “intelligence”.

<sup>2</sup> For an overview, see (Emanuel et al. 1991; Buchanan and Brock 2019).

<sup>3</sup> For discussion, see (Shalowitz, Garrett-Mayer, and Wendler 2006; Jezewski et al. 2007).

In response, researchers have recently suggested an ML based solution.<sup>4</sup> Very roughly, the idea is that we can attempt to accurately model incapacitated patients' preferences; we can then use that model to predict what patients would want under a range of clinical conditions, and then use those predictions as the basis for clinical care. Supposing it is technically feasible to develop the algorithmic part of a patient preference predictor (PPP), such a model would need to be *trained*. Such training is sometimes said to face a serious logistical challenge, viz. somehow (legally, one hopes) acquiring the necessary data.<sup>5</sup> For reasons of space, I'll ignore this challenge in what follows.

Whatever the logistical challenges, there are clear ethical concerns with PPPs.<sup>6</sup> However, extant concerns do not target PPPs *qua* ML models trained using deep learning. Instead, these concerns apply equally well if (say) PPPs are entirely hand-programmed statistical models.<sup>7</sup> In this paper I develop a novel ethical problem for PPPs -- one that applies to PPPs specifically in virtue of their nature as ML models trained using modern deep learning techniques. I'll argue this problem is sufficiently morally serious to shift the normative burden of proof: those in favor of developing and deploying PPPs in a clinical setting owe us a solution before patient-involving experimental research can safely begin. Absent a meaningful risk mitigation measure, institutional review boards will and should ban exactly the practical research needed to move the proposal forward. This puts the moral ball in the ML researcher's technical court: show us a way to assure ourselves against the moral hazard, or risk a halt to progress in an area of potentially important impact for AI in medicine. Here is how I proceed.

Section 2 gives a brief overview of why predictive models trained using deep learning are indifferent as to methods for achieving improvements in accuracy. Here, I highlight the fact that improvements in accuracy can potentially be achieved by so-called "performative prediction." Section 3 lays out the idea of preference shaping and explains why it is morally illicit as a means to achieving accuracy in the case of PPPs. This yields a normative bar that proposals to use PPPs must clear: they must show that a particular model is not incentivized to

---

<sup>4</sup> (Rid and Wendler 2014b; Wendler 2021; Wendler et al. 2016; Brock 2014; Biller-Andorno and Biller 2019; Ferrario, Gloeckler, and Biller-Andorno 2022).

<sup>5</sup> See (Rid and Wendler 2014a) for discussion; for relevant machine learning research, see (O. Evans et al. 2018).

<sup>6</sup> For a selection of moral criticism, see (N. Sharadin 2019; N. P. Sharadin 2018; Ditto and Clark 2014; Kim 2014; John 2014; Dresser 2014; Tretter and Samhammer 2023; Mainz 2022).

<sup>7</sup> Compare (N. P. Sharadin 2018).

improve accuracy by way of preference shaping. Section 4 concludes by connecting the problem identified here to broader concerns about AI safety.

## 2 Reward, Accuracy & Performative Prediction

The performance of predictive models is measured by how close their predictions are to the way the world actually is. For example: if it actually rains 80% of the time a weather forecasting model reports an 80% chance of rain, this is a well-performing model. Call this measure of performance ‘accuracy.’<sup>8</sup> Any *inaccuracy* can in principle be corrected in one of two ways.<sup>9</sup> The first, familiar method, involves changing the model (and so its predictions) to better match the actual distribution of probabilities. This can be done by retraining the model, fine-tuning it, or by otherwise altering its architecture. Less familiarly, but obviously, inaccuracies in a model’s predictions can also be corrected by changing the actual distribution of probabilities, which is to say *changing the world*. For instance, if a climate model predicts that it will be more than 3 degrees warmer by the end of the century, then one way to ensure this prediction is correct is by emitting as much carbon into the atmosphere as humanly possible.

ML models trained using deep learning aim to maximize their expected accuracy (or: to minimize inaccuracy); they are incentivized to be on an entirely accuracy-based metric the best predictor they can be.<sup>10</sup> There are in general no constraints on what counts as the “right kind” of improvement in accuracy: accuracy is accuracy is accuracy, however it’s achieved. One way to think about this feature of deep-learning trained systems is that this kind of training doesn’t typically restrict the *permissible means* to maximizing accuracy.<sup>11</sup> By default, then, models are indifferent as to *how* to maximize accuracy. For instance, if it was possible for a weather model predicting rain to make it rain, or for a climate model predicting warming to heat up the planet, then, in principle at least, *ceteris paribus*, the model would be indifferent to improving the accuracy of its predictions by making it rain or heating up the planet as compared to adjusting its predictions.

---

<sup>8</sup> For a technical overview, see (Gneiting and Raftery 2007).

<sup>9</sup> Well, three. We could change our scoring rule, or our performance metric. I ignore this possibility in what follows.

<sup>10</sup> I follow the literature in saying that a learner is *incentivized* to do something just in case doing that thing increases performance (or reward). See (Krueger, Maharaj, and Leike 2020, 2).

<sup>11</sup> If this sounds familiar from the Forever War between consequentialists and Kantians, that’s not an accident.

This might seem worrying: after all, we can expect models to be indifferent between improvements to accuracy arrived at by changing their *representations* of the probabilities and those arrived at by changing the actual *probabilities* -- the *facts themselves*. But these are two very different kinds of thing, and we certainly don't want our predictive models doing the latter! The natural reaction to this worry is that it's misplaced. There aren't (to our knowledge) any (e.g. weather or climate) models in existence or development that can act on the world in the usual kind of way required to change the actual distribution of probabilities (e.g. by making it rain or emitting carbon). Models don't really *do* anything.

But, natural as it is, this reaction is too quick. Models do at least one thing. They make predictions. And predictions can, after all, affect the world. For instance: a hedge fund's model predicts that NYSE:GME will fall, and as a result the fund publicly takes up a short position (on margin). Maybe this causes other investors to lose confidence, and so causes the stock to fall. Maybe not: it might instead cause part of the internet to lose its collective mind, attempt a squeeze, and so cause the stock to (briefly) go to the moon.<sup>12</sup> Either way, the model didn't just make a prediction, it made a (more or less convoluted) difference to the *actual likelihood that that prediction would be correct*, i.e., it made a difference to the facts on the ground.

Machine learning researchers do not agree on a name for this phenomenon, whereby a model's prediction can affect its own accuracy.<sup>13</sup> Here, I begrudgingly agree to call it "performative prediction."<sup>14</sup> Performative prediction is the family of phenomenon whereby a model's predictions (somehow) make a difference to the spread of probability distributions it aims to represent. Under what conditions are models of the sort we've described here *incentivized* to make performative predictions? That's a trick question; the answer is: under the same conditions they're incentivized to make *any prediction whatsoever*, viz. that doing so maximizes accuracy. They are indifferent between means for maximizing accuracy.

---

<sup>12</sup> See (Good 2021).

<sup>13</sup> Philosophers call a related phenomenon self-fulfilling beliefs (Silva, Paul Jr forthcoming; Antill 2019).

<sup>14</sup> Following (Perdomo et al. 2020). *Begrudgingly* because it can make it sound as if the model *itself* is doing something. It isn't: we are doing something with the model.

### 3 Preference Shaping & Performative Prediction

Preference shaping is when an agent's preferences are induced to change exogenously. Preference shaping *per se* is morally neutral, as in: your preference for cilantro exogenously changes when you move to Mexico City. Of course, sometimes preference shaping is *not* morally neutral, as in: you are beaten daily until you come to enjoy cilantro.<sup>15</sup>

In the clinical context, patients' preferences regarding care are generally regarded as sacrosanct: they ought not be intentionally, wittingly shaped.<sup>16</sup> There are, as everywhere, exceptions. You might try to talk a Christian Scientist into wanting a blood transfusion. The exceptions to the rule are justified by a *trade-off* in values.<sup>17</sup> On the one hand, there are important values (e.g., autonomy) at stake in an agent's preferences being down to her; on the other hand, there are important values (e.g., reducing harm, improving outcomes) at stake in not letting people prefer outcomes that are worse on some objective, non-preference-based measure.

Is it ever morally permissible for a *clinician* to shape a patient's preferences for the sake of improving the accuracy of a predictive model of those preferences? Obviously, no. Consider the following dialogue:

*Doctor:* Great, I have your results. Our in-house-model, Happy Patient, has predicted that you prefer radiation to surgery. Happily, both are equally effective in your case. So, I'll put you down for radiation.

*Patient:* I actually prefer surgery.

*Doctor:* ...<checks notes>...I see...In that case, hmm... how about having a look at these statistics about death during surgery and these gruesome pictures of surgical mishaps.

*Patient:* I'd rather not.

*Doctor:* Look, you're killing Happy Patient's accuracy score, if you could just...

*Patient:* Can I get a referral?

Not only can *clinicians* shape patient preferences, a *model's performative predictions* can also shape patient preferences. This might seem odd and

---

<sup>15</sup> Compare (Franklin et al. 2022).

<sup>16</sup> This follows from broader ideas about the importance of informed consent. For an overview, see (Faden and Beauchamp 1986); for critical discussion, see (Manson and O'Neill 2007).

<sup>17</sup> This is not controversial. See (Li and Chapman 2020) for discussion.

unfamiliar: how could a model's predictions about a patient's preferences count as performative predictions in this sense -- how could they shape a patient's preferences?

In fact this phenomenon isn't odd; at least, it isn't unfamiliar. There are a number of more or less well-studied ways in which people's preferences over outcomes can be shaped by predictions about those preferences, at least when those predictions in some way causally interact with the patients themselves (e.g. by being presented to them). For example, the literature on "nudging" is rife with strategies for affecting people's preferences by (e.g.) presenting them information about prospective choices in a particular order, or by carefully curating the list of alternatives (e.g. by removing or including irrelevant ones).<sup>18</sup> In the present context, information presented to patients could be (e.g.) a model's predictions about what a patient will want; those predictions about preferences can shape the preferences themselves in exactly the way other kinds of (irrelevant) information can shape people's preferences.

Not only is it possible to *imagine* models' predictions about people's preferences shaping those preferences, we know that in fact this happens in the real world. Forget about medical preferences for a moment. Content-recommendation models, such as those that determine the next song, movie, news clip, or other piece of "content" in an algorithmically determined feed, offer a simple illustration of the phenomenon.<sup>19</sup> Here is how it works for content-recommendation models. If you're predicted by a content-recommendation model to like X, then, *ceteris paribus*, you will be shown more X. In a widely recognized phenomenon known as the "mere exposure" effect, "mere" exposure to X is extraordinarily likely to increase your preference for X.<sup>20</sup> Hence, you will (at least on the margins) come to like X, per the model's predictions. Hence, content-recommendation models can (and do!) shape peoples' preferences. And they do so in a way that improves their own predictions: they are *notorious* performative predictors.<sup>21</sup> The present point is

---

<sup>18</sup> For a recent philosophical discussion, see Parmer (2023). The debate over the ethics of nudging is ongoing. For the classic source on "nudges" see Thaler and Sunstein (2008).

<sup>19</sup> For technical discussion of the broad phenomenon, see (Krueger, Maharaj, and Leike 2020; C. Evans and Kasirzadeh 2022; Farquhar, Carey, and Everitt 2022; Everitt et al. 2021).

<sup>20</sup> (Mrkva and Van Boven 2020; Pliner 1982; Zebrowitz, White, and Wieneke 2008; Lakkakula et al. 2010; Pennycook, Cannon, and Rand 2018; Rapp and Salovich 2018; Ulusoy et al. 2021)

<sup>21</sup> See the discussion in (Perdomo et al. 2020).

that PPPs are, in effect, a kind of content-recommendation engine: they recommend *medical content*.<sup>22</sup>

Is it ever morally permissible for the predictions of a PPP to (by *whatever* causal route) shape patient preferences simply in order to improve the accuracy of those very predictions? No. This is for the same reasons as before. The improved accuracy of the predictions doesn't in the requisite way trade-off against the values at stake in shaping patients' preferences.

This doesn't mean that, as a moral matter, PPPs must never actually performatively predict in a way that turns out to shape patients' preferences (and so perhaps improve accuracy). That would be too high a bar, as the content recommendation example, together with the research on the mere exposure effect, illustrates. Simply being told that you prefer something is itself somewhat likely to make you prefer it. And if PPPs are to be *used* (rather than stuck in a drawer) then their predictions will presumably have *some* causal impact, that causal impact might involve shaping patients' preferences, and it might thereby improve the accuracy of the PPP.<sup>23</sup>

Equally: a clinician might *actually* shape a patient's preferences and they might do so in a way that improves a predictive model's accuracy. What's morally impermissible is shaping a patient's preferences *as a means to improving the accuracy of a predictive model*. To avoid this moral hazard in the case of a clinician, we simply detach any incentives a clinician might have for improving the accuracy of a model from the incentives they have to shape patients' preferences. In effect, we disallow or disincentivize dialogues like the one above. There are many obvious ways to do this.

How do we avoid this moral hazard in the case of a model trained using deep learning? That's a very good question. It is trivial to describe the property we want a suitable model to have. We want it to be such that it is manifestly, provably not incentivized to make performative predictions that shape patient preferences. But it turns out to be extraordinarily difficult to assure ourselves that any given model in fact has this property. There are, for example, no extant proposals in the machine learning literature for identifying the conditions

---

<sup>22</sup> Thanks to an anonymous referee for encouraging clarity on this point.

<sup>23</sup> The only research that I'm aware of that approaches the question of performative prediction in the context of medical AI is a review article (Chen et al. 2021); there, the authors simply note the possibility of distributional shift (aka performative prediction).

under which a model has this property.<sup>24</sup> This is a Bad Thing. If we lack a meaningful way to assure ourselves that a particular PPP is not incentivized to shape patients' preferences, then it seems clear that experimental, which is to say patient-involving, research on PPPs will (and should) be blocked by institutional review boards, which correctly take a dim view of this kind of moral risk. It goes without saying that PPPs should not be deployed in a clinical setting.

#### 4 Conclusion & Discussion

Let me step back from the particulars for a moment; below, I'll return to them. Readers familiar with the broader literature on AI safety will not be surprised by anything they've read. There's a concatenation of long-standing, well-known, very hard problems in AI safety that all have something like the general form of the problem I've here identified in a particular case: either we can't actually affect a model's incentives, or we can't interrogate them, or having done either of those we can't *assure* ourselves of the precise content of those incentives, etc. Despite important differences in theorizing about and technical approaches to solving these and related problems, they're often lumped under one name: the "Alignment Problem".<sup>25</sup> The alignment problem is big, fuzzy, and poorly understood. So, one way to respond to what I've said so far is to point out that it's simply an instance of a well-known (if not well-understood) problem and, moreover, *give us a minute, we're working on it*.

That reaction, I think, is a mistake given the present context. The alignment problem may be big and fuzzy, but the problem identified here is relatively small and precise.<sup>26</sup> ML models are being used in medicine right now, today. PPPs are being proposed not just as an interesting idea, but as a thing that should begin to be put into practice.<sup>27</sup> We should not wait for a solution to the broadest possible description of the broadest possible AI safety problem (e.g., the Alignment Problem) to clearly articulate the moral hazards involved in particular proposed uses of the technology. I am not sure I know what progress

---

<sup>24</sup> This is also the conclusion of other AI safety researchers. Compare (Hendrycks et al. 2022; C. Evans and Kasirzadeh 2022; Ashton and Franklin 2022). This is not to say that there are no proposals about how to ensure that models have *other* interesting properties related to performative prediction, e.g., can achieve various strategic equilibria; see (Mendler-Dünner et al. 2020; Brown, Hod, and Kalemaj 2022; Miller, Perdomo, and Zrnic 2021).

<sup>25</sup> This is part of why no one agrees on a definition of *the* Alignment Problem.

<sup>26</sup> Thanks to an anonymous referee for this way of putting the contrast between the alignment problem and the problem I identify in the paper.

<sup>27</sup> (Wendler 2021; Ferrario, Gloeckler, and Biller-Andorno 2022).



on the Alignment Problem looks like, or even what way it is best to talk about many of the questions that researchers seem to care about in this broad area. But things are much simpler in this case. A proposal is being seriously floated in the scientific and philosophical literature to deploy technology that will, by the technical experts' own admission, not be disincentivized from bringing about what is an uncontroversially serious moral harm.<sup>28</sup>

## References

- Antill, Gregory. 2019. "Evidence and Self-Fulfilling Belief." *American Philosophical Quarterly* 56 (4): 319–30. <https://doi.org/10.2307/48563046>.
- Biller-Andorno, Nikola, and Armin Biller. 2019. "Algorithm-Aided Prediction of Patient Preferences - An Ethics Sneak Peek." *The New England Journal of Medicine* 381 (15): 1480–85. <https://doi.org/10.1056/NEJMms1904869>.
- Brock, Dan W. 2014. "Reflections on the Patient Preference Predictor Proposal." *Journal of Medicine and Philosophy* 39 (2): 153–60.
- Brown, Gavin, Shlomi Hod, and Iden Kalemaj. 2022. "Performative Prediction in a Stateful World." In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, 6045–61. PMLR. <https://proceedings.mlr.press/v151/brown22a.html>.
- Buchanan, Allen, and Dan W. Brock. 2019. "Deciding for Others." *Death, Dying and the Ending of Life, Volumes I and II*, 205–82.
- Carroll, Micah D., Anca Dragan, Stuart Russell, and Dylan Hadfield-Menell. 2022. "Estimating and Penalizing Induced Preference Shifts in Recommender Systems." In *Proceedings of the 39th International Conference on Machine Learning*, 2686–2708. PMLR. <https://proceedings.mlr.press/v162/carroll22a.html>.
- Chen, Irene Y., Shalmali Joshi, Marzyeh Ghassemi, and Rajesh Ranganath. 2021. "Probabilistic Machine Learning for Healthcare." *Annual Reviews of Biomedical Data Science*. <https://doi.org/10.48550/arXiv.2009.11087>.
- Ditto, Peter H., and Cory J. Clark. 2014. "Predicting End-of-Life Treatment Preferences: Perils and Practicalities." *The Journal of Medicine and Philosophy: A Forum for Bioethics and Philosophy of Medicine* 39 (2): 196–204. <https://doi.org/10.1093/jmp/jhu007>.
- Dresser, Rebecca. 2014a. "Law, Ethics, and the Patient Preference Predictor." *Journal of Medicine and Philosophy* 39 (2): 178–86.

---

<sup>28</sup> Thanks to Simon Goldstein, Dan Hendrycks, Jacqueline Harding, Cameron Kirk-Giannini, David Krueger, Nick Laskowski, Robert Long, Elliot Thornley, and members of the Cottage Group for helpful discussions about these and related issues.

- . 2014b. “Law, Ethics, and the Patient Preference Predictor.” *The Journal of Medicine and Philosophy* 39 (2): 178–86. <https://doi.org/10.1093/jmp/jhu004>.
- Emanuel, Linda L., Michael J. Barry, John D. Stoeckle, Lucy M. Ettelson, and Ezekiel J. Emanuel. 1991. “Advance Directives for Medical Care — A Case for Greater Use.” *New England Journal of Medicine* 324 (13): 889–95. <https://doi.org/10.1056/NEJM199103283241305>.
- Evans, Charles, and Atoosa Kasirzadeh. 2022. “User Tampering in Reinforcement Learning Recommender Systems.” In . arXiv. <https://doi.org/10.48550/arXiv.2109.04083>.
- Evans, Owain, Andreas Stuhlmüller, Chris Cundy, Ryan Carey, Thomas McGrath, and Andrew Schreiber. 2018. “Predicting Human Deliberative Judgments with Machine Learning.”
- Everitt, Tom, Marcus Hutter, Ramana Kumar, and Victoria Krakovna. 2021. “Reward Tampering Problems and Solutions in Reinforcement Learning: A Causal Influence Diagram Perspective.” *Synthese* 198 (27): 6435–67. <https://doi.org/10.1007/s11229-021-03141-4>.
- Faden, Ruth R., and Tom L. Beauchamp. 1986. *A History and Theory of Informed Consent*. Oxford University Press.
- Farquhar, Sebastian, Ryan Carey, and Tom Everitt. 2022. “Path-Specific Objectives for Safer Agent Incentives.” arXiv. <https://doi.org/10.48550/arXiv.2204.10018>.
- Fazio, Lisa K., Nadia M. Brashier, B. Keith Payne, and Elizabeth J. Marsh. 2015. “Knowledge Does Not Protect against Illusory Truth.” *Journal of Experimental Psychology: General* 144: 993–1002. <https://doi.org/10.1037/xge0000098>.
- Ferrario, Andrea, Sophie Gloeckler, and Nikola Biller-Andorno. 2022. “Ethics of the Algorithmic Prediction of Goal of Care Preferences: From Theory to Practice.” *Journal of Medical Ethics*, November, medethics-2022-108371. <https://doi.org/10.1136/jme-2022-108371>.
- Franklin, Matija, Hal Ashton, Rebecca Gorman, and Stuart Armstrong. 2022. “Recognising the Importance of Preference Change: A Call for a Coordinated Multidisciplinary Research Effort in the Age of AI.” *The AAAI-22 Workshop on AI For Behavior Change*, 1–7. <https://doi.org/10.48550/arXiv.2203.10525>.
- Gneiting, Tilmann, and Adrian E Raftery. 2007. “Strictly Proper Scoring Rules, Prediction, and Estimation.” *Journal of the American Statistical Association* 102 (477): 359–78. <https://doi.org/10.1198/016214506000001437>.

- Good, Owen S. 2021. "GameStop's Stock Market Explosion, Explained." Polygon (blog). January 27, 2021. <https://www.polygon.com/2021/1/27/22252600/gamestop-stock-gme-why-whats-happening-explain>.
- Healy, Kieran. 2015. "The Performativity of Networks." *European Journal of Sociology / Archives Européennes de Sociologie* 56 (2): 175–205. <https://doi.org/10.1017/S0003975615000107>.
- Hendrycks, Dan, Nicholas Carlini, John Schulman, and Jacob Steinhardt. 2022. "Unsolved Problems in ML Safety." arXiv. <https://doi.org/10.48550/arXiv.2109.13916>.
- Jardas, E. J., David Wasserman, and David Wendler. 2022. "Autonomy-Based Criticisms of the Patient Preference Predictor." *Journal of Medical Ethics* 48 (5): 304–10.
- Jezewski, Mary Ann, Mary Ann Meeker, Loralee Sessanna, and Deborah S. Finnell. 2007. "The Effectiveness of Interventions to Increase Advance Directive Completion Rates." *Journal of Aging and Health* 19 (3): 519–36. <https://doi.org/10.1177/0898264307300198>.
- John, Stephen. 2014. "Patient Preference Predictors, Apt Categorization, and Respect for Autonomy." *The Journal of Medicine and Philosophy: A Forum for Bioethics and Philosophy of Medicine* 39 (2): 169–77. <https://doi.org/10.1093/jmp/jhu008>.
- Kim, Scott Y. H. 2014. "Improving Medical Decisions for Incapacitated Persons: Does Focusing on 'Accurate Predictions' Lead to an Inaccurate Picture?" *The Journal of Medicine and Philosophy: A Forum for Bioethics and Philosophy of Medicine* 39 (2): 187–95. <https://doi.org/10.1093/jmp/jhu010>.
- Krueger, David, Tegan Maharaj, and Jan Leike. 2020. "Hidden Incentives for Auto-Induced Distributional Shift." arXiv. <https://doi.org/10.48550/arXiv.2009.09153>.
- Lakkakula, Anantha, James Geaghan, Michael Zanovec, Sarah Pierce, and Georgianna Tuuri. 2010. "Repeated Taste Exposure Increases Liking for Vegetables by Low-Income Elementary School Children." *Appetite* 55 (2): 226–31. <https://doi.org/10.1016/j.appet.2010.06.003>.
- Li, Meng, and Gretchen B. Chapman. 2020. "Medical Decision Making." In *The Wiley Encyclopedia of Health Psychology*, 347–53. John Wiley & Sons, Ltd. <https://doi.org/10.1002/9781119057840.ch84>.
- Mainz, Jakob Thrane. 2022a. "The Patient Preference Predictor and the Objection from Higher-Order Preferences." *Journal of Medical Ethics*.

- . 2022b. “The Patient Preference Predictor and the Objection from Higher-Order Preferences.” *Journal of Medical Ethics*, July. <https://doi.org/10.1136/jme-2022-108427>.
- Manson, Neil C., and O. O’Neill. 2007. *Rethinking Informed Consent in Bioethics*. Cambridge University Press. <https://eprints.lancs.ac.uk/id/eprint/4042/>.
- Mendler-Dünner, Celestine, Juan Perdomo, Tijana Zrnic, and Moritz Hardt. 2020. “Stochastic Optimization for Performative Prediction.” In *Advances in Neural Information Processing Systems*, 33:4929–39. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2020/hash/33e75ff09dd601bbe69f351039152189-Abstract.html>.
- Mertens, Stephanie, Mario Herberz, Ulf J. J. Hahnel, and Tobias Brosch. 2022. “The Effectiveness of Nudging: A Meta-Analysis of Choice Architecture Interventions across Behavioral Domains.” *Proceedings of the National Academy of Sciences of the United States of America* 119 (1): e2107346118. <https://doi.org/10.1073/pnas.2107346118>.
- Miller, John P., Juan C. Perdomo, and Tijana Zrnic. 2021. “Outside the Echo Chamber: Optimizing the Performative Risk.” In *Proceedings of the 38th International Conference on Machine Learning*, 7710–20. PMLR. <https://proceedings.mlr.press/v139/miller21a.html>.
- Mrkva, Kellen, and Leaf Van Boven. 2020. “Salience Theory of Mere Exposure: Relative Exposure Increases Liking, Extremity, and Emotional Intensity.” *Journal of Personality and Social Psychology* 118: 1118–45. <https://doi.org/10.1037/pspa0000184>.
- Parmer, W. Jared (2023). Nudges, Nudging, and Self-Guidance Under the Influence. *Ergo* 9 (44):1199-1232.
- Pennycook, Gordon, Tyrone D. Cannon, and David G. Rand. 2018. “Prior Exposure Increases Perceived Accuracy of Fake News.” *Journal of Experimental Psychology. General* 147 (12): 1865–80. <https://doi.org/10.1037/xge0000465>.
- Perdomo, Juan, Tijana Zrnic, Celestine Mendler-Dünner, and Moritz Hardt. 2020. “Performative Prediction.” In *Proceedings of the 37th International Conference on Machine Learning*, 7599–7609. PMLR. <https://proceedings.mlr.press/v119/perdomo20a.html>.
- Pliner, Patricia. 1982. “The Effects of Mere Exposure on Liking for Edible Substances.” *Appetite* 3 (3): 283–90. [https://doi.org/10.1016/S0195-6663\(82\)80026-3](https://doi.org/10.1016/S0195-6663(82)80026-3).
- Rapp, David N., and Nikita A. Salovich. 2018. “Can’t We Just Disregard Fake News? The Consequences of Exposure to Inaccurate Information.” *Policy*

- Insights from the Behavioral and Brain Sciences* 5 (2): 232–39.  
<https://doi.org/10.1177/2372732218785193>.
- Rid, Annette, and David Wendler. 2014a. “Treatment Decision Making for Incapacitated Patients: Is Development and Use of a Patient Preference Predictor Feasible?” *Journal of Medicine and Philosophy* 39 (2): 130–52.
- . 2014b. “Use of a Patient Preference Predictor to Help Make Medical Decisions for Incapacitated Patients.” *Journal of Medicine and Philosophy* 39 (2): 104–29.
- Salmond, Susan W., and Estrella David. 2005. “Attitudes Toward Advance Directives and Advance Directive Completion Rates.” *Orthopaedic Nursing* 24 (2): 117.
- Shalowitz, David I., Elizabeth Garrett-Mayer, and David Wendler. 2006. “The Accuracy of Surrogate Decision Makers: A Systematic Review.” *Archives of Internal Medicine* 166 (5): 493–97.
- Sharadin, Nathaniel. 2019. “Should Aggregate Patient Preference Data Be Used to Make Decisions on Behalf of Unrepresented Patients?” *AMA Journal of Ethics* 21 (7): 566–74.
- Sharadin, Nathaniel Paul. 2018. “Patient Preference Predictors and the Problem of Naked Statistical Evidence.” *Journal of Medical Ethics* 44 (12): 857–62.  
<https://doi.org/10.1136/medethics-2017-104509>.
- Silva, Paul Jr, Paul Silva. forthcoming. “Self-Fulfilling Beliefs: A Defense.” *Australasian Journal of Philosophy*.
- Tandoc Jr., Edson C. 2019. “The Facts of Fake News: A Research Review.” *Sociology Compass* 13 (9): e12724. <https://doi.org/10.1111/soc4.12724>.
- Thaler, Richard H., Sunstein, Cass R.. *Nudge: Improving Decisions about Health, Wealth and Happiness*. United Kingdom: Penguin Books, 2008.
- Tretter, Max, and David Samhammer. 2023. “For the Sake of Multifacetedness. Why Artificial Intelligence Patient Preference Prediction Systems Shouldn’t Be for next of Kin.” *Journal of Medical Ethics*, January.  
<https://doi.org/10.1136/jme-2022-108775>.
- Ulusoy, Ezgi, Dustin Carnahan, Daniel E. Bergan, Rachel C. Barry, Siyuan Ma, Suhwoo Ahn, and Johnny McGraw. 2021. “Flooding the Zone: How Exposure to Implausible Statements Shapes Subsequent Belief Judgments.” *International Journal of Public Opinion Research* 33 (4): 856–72.  
<https://doi.org/10.1093/ijpor/edab022>.
- Wendler, David. 2021. “A Call for a Patient Preference Predictor.” *Critical Care Medicine* 49 (6): 877–80.
- Wendler, David, Bob Wesley, Mark Pavlick, and Annette Rid. 2016. “A New Method for Making Treatment Decisions for Incapacitated Patients: What

Do Patients Think about the Use of a Patient Preference Predictor?" *Journal of Medical Ethics* 42 (4): 235–41.

Zebrowitz, Leslie A., Benjamin White, and Kristin Wieneke. 2008. "Mere Exposure and Racial Prejudice: Exposure to Other-Race Faces Increases Liking for Strangers of That Race." *Social Cognition* 26 (3): 259–75. <https://doi.org/10.1521/soco.2008.26.3.259>.