

1 We thank the reviewers for their excellent feedback and for recognizing the value of our work to the NeurIPS audience.

2 **R1, R2: Realism of the causal graph.** In our model, the outcome Y_t is determined not only by A_t , but also by the
3 unobserved moderating variable Z and possibly exogenous variables. Z may capture, for example, the disease state
4 of a subject and introduce *correlations* between outcomes Y_s and Y_t for $s < t$. It is a common working assumption
5 that previous actions and outcomes, A_s, Y_s for $s < t$, do not have direct *causal* effect on Y_t when Y_t represents the
6 symptoms of chronic conditions where drugs do not affect the underlying disease state; given sufficient time between
7 treatments, symptoms return to a baseline level until another treatment is started. In addition to rheumatoid arthritis [1],
8 which is used as motivation in the manuscript, examples include depression [2] and Parkinson’s disease [3].

9 **R2: Usefulness of causal framework.** The causal framework is necessary for distinguishing the observational
10 outcome Y_t from the potential outcome $Y(a)$. Only under certain assumptions can $Y(a)$ be estimated from Y_t .
11 Sufficient conditions for this are given in Section 4. The distinction is particularly important in the case where $Y(a)$ is
12 not fully identifiable from observational data due to unobserved confounding, as discussed in Section 4.2.

13 **R2: Observed and counterfactual outcomes.** The relationship between counterfactual outcomes $Y(a)$ and observed
14 outcomes Y_t is established in Theorem 1. The LHS of (4), $\rho(h_s)$, is a function of counterfactual outcomes, as defined
15 in (2), and the RHS is a function of observed outcomes Y_t . As remarked after Theorem 1, under our assumptions, a
16 model of $p(Y_t | H_t = h, A_t = a)$ is sufficient for estimating the distribution of $Y(a)$. To solve our policy optimization
17 problem (1), it is not necessary to impute all counterfactuals. For example, in the binary case given by R2, if an action
18 $a = 0$ has been tried, and $Y(0)$ observed, only the probability that $Y(1) > Y(0)$ is required to solve the problem.

19 **R2: Relation to model-free RL.** R2 is correct that the model-free method NDP compares similarly to the other
20 methods in the antibiotics experiment. However, we show in Figure 1b that the qualitative behavior of NDP as a function
21 of dataset size is very different from that of CDP and CG. Additionally, as shown in Appendix A.7 (Thm A7), NDP is
22 suboptimal in the general case. We hope that these contributions are recognized. By a “transparent” tradeoff, we refer to
23 the meaning of the parameters δ (CDP, CG) and λ (NDP). δ is directly interpretable as a probability threshold at which
24 we are satisfied with the best-so-far treatment (as used in the antibiotics experiment). The value of λ does not have an
25 immediate interpretation as a level of certainty of near-optimality—the *tradeoff for a fixed λ varies across datasets*.

26 **R3, R4: Comparison with experts and the emulated expert.** R3 is correct that it is feasible in practice to include
27 more information in the policy so that it compares more favorably to experts in the first step. Similarly, the accuracy of
28 the emulated doctor, remarked upon by R4, could be improved by using more information to emulate the doctor policy.
29 In our experiment, we intentionally kept the patient representation small because the number of samples was fairly
30 limited. The emulated expert in our study attempts to approximate the expert’s policy with the same information given
31 to the other algorithms. As such, it serves as an *imitation learning* baseline to complement the policy optimization
32 approaches developed in this work. We will clarify this choice in the paper. We note, however, that our algorithm is
33 trying to achieve a different goal than the expert. While experts may attempt to prescribe the best action on the first try,
34 we try to minimize the expected number of tested treatments. In that sense, the expert can be thought of as a greedy
35 agent, which our paper argues is not always optimal, if we’re trying to minimize needless trial and error on patients.

36 **R4: Impact of this work.** We certainly agree that our method is not suitable for all applications. However, there is a
37 large class of medical conditions and treatments which fall exactly under the specifications of our model. In fact, our
38 motivation for this work is the result of working with active clinicians treating rheumatoid arthritis. The application
39 of our method for this purpose is an ongoing project which will be aimed at the clinical research community, but we
40 appreciate the reviewer’s feedback which highlighted that the use cases of our method are not readily apparent from
41 reading the paper. We will add a description of the RA problem to the paper to make it more concrete and demonstrate
42 a motivating use case, along with a discussion of other uses such as in treating psychiatric disorders. Finally, we believe
43 that the problem has applications also outside of medicine, such as for general recommendation systems.

44 **R4: Short-term response.** It is true that short-term response is critical for some applications and should not be
45 discounted; this is a potential challenge also for reinforcement learning which optimizes long-term return, possibly
46 sacrificing immediate rewards. Our goal is to find a near-optimal treatment in as few steps as possible which is an
47 important consideration in other applications [2]. If an optimal treatment *can* be reliably identified in a single step,
48 the algorithm is incentivized to do so. Short-term success is sacrificed only if there is great uncertainty about which
49 treatment is likely to work and this can be reduced by a sub-optimal treatment. Our greedy approximation incentivizes
50 short-term response by preferring actions higher that are likely to have a higher outcome (see 5.2). Such incentives
51 could be incorporated into the dynamic programming solution as well, and is an interesting direction for future work.

52 [1] J. R. O’Dell. Therapeutic strategies for rheumatoid arthritis. *New England Journal of Medicine*, 350(25):2591–2602, 2004.

53 [2] M. F. Pradier, T. H. McCoy Jr, M. Hughes, R. H. Perlis, and F. Doshi-Velez. Predicting treatment dropout after antidepressant initiation. *Translational psychiatry*,
54 10(1):1–8, 2020.

55 [3] M. Stacy, A. Bowron, M. Guttman, R. Hauser, K. Hughes, J. P. Larsen, P. LeWitt, W. Oertel, N. Quinn, K. Sethi, et al. Identification of motor and nonmotor
56 wearing-off in parkinson’s disease: comparison of a patient questionnaire versus a clinician assessment. *Movement disorders: official journal of the Movement*
57 *Disorder Society*, 20(6):726–733, 2005.