
Cognitive Information Filters: Algorithmic Choice Architecture for Boundedly Rational Choosers

Stefan Bucher^{1,2} Peter Dayan^{1,2}

¹Max Planck Institute for Biological Cybernetics ²University of Tübingen AI Center
72076 Tübingen, Germany
{stefan.bucher, peter.dayan}@tuebingen.mpg.de

Abstract

We introduce *cognitive information filters* as an algorithmic approach to mitigating information overload using choice architecture: We develop a rational inattention model of boundedly rational multi-attribute choice and leverage it to programmatically select information that is effective in inducing desirable behavioral outcomes. By inferring preferences and cognitive constraints from boundedly rational behavior, our methodology can optimize for revealed welfare and hence promises better alignment with boundedly rational users than recommender systems optimizing for imperfect welfare proxies such as engagement. This has implications beyond economics, for example for alignment research in artificial intelligence.

1 Introduction

Information overload is ubiquitous, as constraints on our cognitive capacity to process information adversely impact the quality of the decisions we make. Limited attention and cognition have important consequences for welfare and markets [McFadden, 2023], and have consequently become of central interest to economists and policy-makers concerned that “individuals are able to pay only limited attention to important aspects of their environment, often have a difficult time processing information, and make cognitive errors even in simple situations”, as stated in the National Academies’ recent consensus study report on behavioral economics [Buttenheim et al., 2023, p. 7].

Behavioral scientists often rely on nudges such as choice architecture to aid boundedly rational choosers in making better decisions, but identifying reliable nudges with significant effect sizes can be difficult and costly [DellaVigna and Linos, 2022, Mertens et al., 2022, Maier et al., 2022] – particularly with heterogeneous populations. In contrast, algorithmic recommender systems are designed to learn from observing individual user choices, but they often suffer from misalignment due to the difficulty of inferring users’ preferences from their boundedly rational behavior [Kleinberg et al., 2022].

In this paper, we introduce *cognitive information filters* as a principled, algorithmic approach to mitigating information overload, based on an information-theoretic model of decision-making under cognitive costs. Specifically, we first develop a rational inattention [Sims, 2003, Maćkowiak et al., 2023] model of multi-attribute choice to describe the behavior of a decision-maker (receiver) facing cognitive information processing costs. We then use model-based, online reinforcement learning to solve the information design problem of a sender choosing which options and features (attributes) are accessible to the rationally inattentive receiver, in order to nudge or persuade them. Observing only the receiver’s choices, the sender learns from repeated interactions which information is most effective in attaining desirable receiver choices, by inferring the receiver’s unobservable preferences and cognitive constraints.

Our approach explicitly addresses the challenge of inferring the preferences of boundedly rational users, by building on recent advances in cognitive economics [e.g. [Woodford, 2020](#), [Caplin, 2023](#)] to develop a more principled approach to choice architecture and information design. Careful information-theoretic modeling opens the door for algorithmically tailoring information to decision-makers’ revealed preferences and cognitive constraints. By allowing the sender to optimize for revealed welfare, this approach promises to be (1) less paternalistic than traditional nudging techniques, and (2) less susceptible to misalignment than recommender systems that optimize for imperfect welfare proxies such as engagement.

2 Model

We model the repeated interaction of a sender (“Alice”) with a receiver (“Bob”). In each period t , Bob makes a choice a_t from a grand set of m choice options with n features characterized by a matrix $\mathbf{X}_t \in \mathcal{X} := \mathbb{R}^{m \times n}$. Bob’s utility $u_w(a, \mathbf{X}) = \mathbf{e}_a^T \mathbf{X} \mathbf{w}$ is linear with preference weights $\mathbf{w} \in \mathbb{R}^n$ and \mathbf{e}_a the a -th standard basis vector. Bob does not observe \mathbf{X}_t , but is rationally inattentive with a prior belief $\mathbf{X}_t \stackrel{iid}{\sim} \mu \in \Delta(\mathbb{R}^{m \times n})$ and marginal information costs κ . Bob can acquire costly information on a subset $\mathbf{Y}_t \in \mathcal{Y} := \mathbb{R}^{k \times l}$ of $k \leq m$ options and $l \leq n$ features.

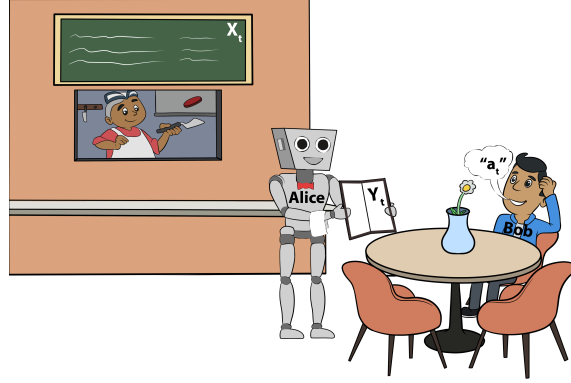


Figure 1: Illustration of the model. In each period, Alice chooses the subset \mathbf{Y}_t of \mathbf{X}_t on which Bob can acquire costly information before making a choice a_t .

This subset is determined by Alice’s choice of an *action selection matrix* $\mathbf{A}_t \in \mathbb{R}^{k \times m}$ and a *feature selection matrix* $\mathbf{F}_t \in \mathbb{R}^{l \times n}$ as $\mathbf{Y}_t = \mathbf{A}_t \mathbf{X}_t \mathbf{F}_t^T$, which we refer to as the *filtered state*. Before choosing the filters \mathbf{A}_t and \mathbf{F}_t , Alice observes \mathbf{X}_t , but not Bob’s type $\theta := (\mathbf{w}, \kappa) \in \Theta$, which is static and distributed according to $\theta \sim \tau \in \Delta(\Theta)$. Alice maintains a belief $b_t \in \Delta(\Theta)$ about Bob’s type, with $b_0 = \tau$; her utility $v(a_t, \mathbf{X}_t, \theta)$ may or may not be aligned with Bob’s.

Given \mathbf{A}_t and \mathbf{F}_t , Bob (who is strategically naïve with respect to Alice) chooses a state-dependent stochastic choice function $P_t : \mathcal{Y} \rightarrow \Delta(A(\mathbf{A}_t))$, i.e. a distribution over the available choice set $A(\mathbf{A}_t)$ conditional on the filtered state \mathbf{Y}_t , so as to maximize expected utility net of information costs $K_\theta(P; \mathbf{A}, \mathbf{F}, \mu) = \kappa I_P(a; \mathbf{Y})$, which are linear in the Shannon mutual information between the filtered state and action.¹ Bob’s problem is thus

$$P_t = \arg \max_{P_t} \int_{\mathcal{X}} \sum_{a_t \in A(\mathbf{A}_t)} P_t(a_t | \mathbf{A}_t \mathbf{X}_t \mathbf{F}_t^T; \theta, \mathbf{A}_t, \mathbf{F}_t, \mu) u_\theta(a_t, \mathbf{A}_t \mathbf{X}_t) d\mu(\mathbf{X}_t) - K_\theta(P_t; \mathbf{A}_t, \mathbf{F}_t, \mu) \quad (1)$$

The timing in each period $t \in \{0, \dots, T\}$ is as follows:

1. Alice observes the realization of $\mathbf{X}_t \sim \mu$, but not θ (instead maintaining belief $b_t \in \Delta(\Theta)$).
2. Given \mathbf{X}_t and b_t , Alice chooses \mathbf{A}_t and \mathbf{F}_t .
3. Bob observes θ , \mathbf{A}_t , and \mathbf{F}_t , but not \mathbf{X}_t (instead maintaining belief $\mu \in \Delta(\mathcal{X})$).
4. Bob chooses $P_t(\cdot | \mathbf{Y}_t; \theta, \mathbf{A}_t, \mathbf{F}_t, \mu)$.
5. $a_t \sim P_t(\cdot | \mathbf{Y}_t; \theta, \mathbf{A}_t, \mathbf{F}_t, \mu)$ is realized, Bob receives $u_\theta(a_t, \mathbf{A}_t \mathbf{X}_t) - K_\theta(P; \mathbf{A}_t, \mathbf{F}_t, \mu)$. Alice observes a_t , forms posterior belief $b_{t+1}(\theta | a_0, \dots, a_t)$, and receives utility $\int_{\Theta} b_{t+1}(\theta) [v(a_t, \mathbf{A}_t \mathbf{X}_t, \theta) - (1 - \alpha) K_\theta(P_t; \mathbf{A}_t, \mathbf{F}_t, \mu)] d\theta$, internalizing Bob’s information costs with a discount factor $\alpha \in [0, 1]$.

¹This formulation is well-known to be equivalent to the choice of a costly Blackwell information structure and making a choice contingent on the signal realization [e.g. [Matejka and McKay, 2015](#), Corollary 1].

3 Results

Solving Bob’s Problem Bob’s problem is a non-trivial rational inattention problem with a matrix state and the added complexity that the information structure is restricted to conditioning on the accessible part \mathbf{Y}_t of the state \mathbf{X}_t . The following novel result generalizes [Bucher and Caplin \[2021\]](#) to provide an interpretable solution in analytical closed form. The only condition we impose on the prior is that it is row-exchangeable (i.e., across actions): for any permutation matrix \mathbf{P} it must be the case that $\mu(\mathbf{X}) = \mu(\mathbf{P}\mathbf{X})$ for all $\mathbf{X} \in \mathcal{X}$. This assumption allows for arbitrary statistical dependency across features, and for values to be correlated across actions as long as they are exchangeable.

Theorem 1 *Given a row-exchangeable prior μ , Bob’s problem (eq. 1) has a solution*

$$P_t(a_t | \mathbf{Y}_t; \theta, \mathbf{A}_t, \mathbf{F}_t, \mu) = \frac{z_\theta(a_t, \mathbf{Y}_t; \mathbf{A}_t, \mathbf{F}_t, \mu)}{\sum_{c \in A(\mathbf{A}_t)} z_\theta(c, \mathbf{Y}_t; \mathbf{A}_t, \mathbf{F}_t, \mu)} \quad \delta_{a_t} \in A(\mathbf{A}_t)$$

and $P_t(a_t | \mathbf{Y}_t; \theta, \mathbf{A}_t, \mathbf{F}_t, \mu) = 0$ for all $a_t \notin A(\mathbf{A}_t)$, where

$$z_\theta(a_t, \mathbf{Y}_t; \mathbf{A}_t, \mathbf{F}_t, \mu) = \exp\left(\frac{1}{\kappa} \mathbf{e}_{a_t}^T \mathbf{A}_t \mathbb{E}_{\mu_{\mathbf{X}|\mathbf{Y}}}[\mathbf{X}_t | \mathbf{Y}_t, \mathbf{A}_t, \mathbf{F}_t] \mathbf{w}\right).$$

Note in particular that Bob’s choice behavior depends on his posterior belief about \mathbf{X}_t after observing the realization of a costly signal on \mathbf{Y}_t , so Bob can make an inference about hidden features from observing accessible ones, as long as they are correlated. For a proof, we refer to [Bucher and Dayan](#).

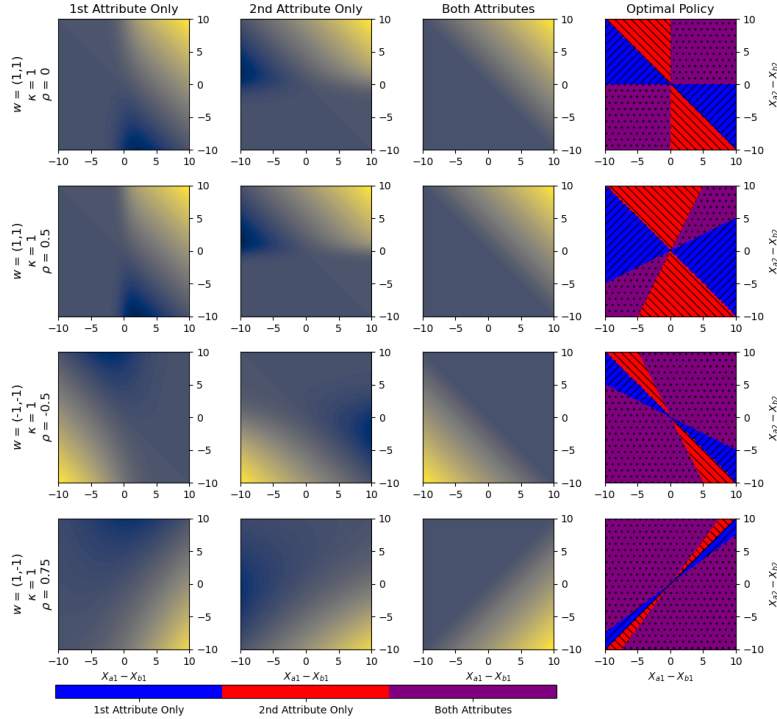


Figure 2: Bob’s gross welfare under three different information policies \mathbf{F} (columns 1-3) along with the welfare-maximizing information policy (column 4), for an example with two options and two features, as a function of $X_{a1} - X_{b1}$ and $X_{a2} - X_{b2}$. Each row corresponds to a set of parameters including Bob’s preference \mathbf{w} and the correlation ρ across features under his prior μ .

Figure 2 illustrates, for an example with two options and two features, how the welfare resulting from Bob’s choices depends on his preference weights w and his prior belief μ (across the figure’s rows) as well as on Alice’s choice of \mathbf{F}_t (across columns). Note how our model captures information overload: The rationally inattentive Bob may, depending on the realization of \mathbf{X}_t , be better off (in terms of gross welfare) when one of the attributes is occluded. When choosing \mathbf{A}_t and \mathbf{F}_t , Alice does not know Bob’s type, however.

Solving Alice’s Problem Alice faces the dynamic problem (with discount factor γ) of learning which filters result in desirable behavior while inferring Bob’s type from his behavior, a form of inverse reinforcement learning. We cast Alice’s problem as a partially observable Markov decision process (POMDP), which gives rise to the Bellman optimality equation

$$V(\mathbf{X}, b, \mu) = \max_{\mathbf{A}, \mathbf{F}} E_{\theta, b} \left[\sum_{a \in \mathcal{A}(\mathbf{A})} P(a | \mathbf{A} \mathbf{X} \mathbf{F}^T; \theta, \mathbf{A}, \mathbf{F}, \mu) (R_{\alpha}(a, \mathbf{X}, \theta, \mathbf{A}, \mathbf{F}; \mu) + \gamma E_{\mathbf{X}^0} \mu [V(\mathbf{X}^0, b^0(a), \mu)]) \right]$$

where $R_{\alpha}(a, \mathbf{X}, \theta, \mathbf{A}, \mathbf{F}; \mu) = v(a, \mathbf{A} \mathbf{X}, \theta) - (1 - \alpha) K_{\theta}(P; \mathbf{A}, \mathbf{F}, \mu)$ and $b^0(a)$ is the posterior belief upon observing Bob’s choice of a .

Trading off exploration and exploitation, Alice should experiment efficiently with different information designs in order to learn in repeated interactions what information is most effective in light of Bob’s latent preferences and information costs.

To solve Alice’s problem we rely on an online reinforcement learning algorithm based on Monte Carlo planning [Silver and Veness, 2010]. Figure 3 shows, for an example with binary Θ , simulated sample trajectories for the cumulative mean regret under Alice’s information policy compared to the full-information policy as a benchmark, demonstrating how Alice learns to achieve a higher mean reward than the full-information baseline.

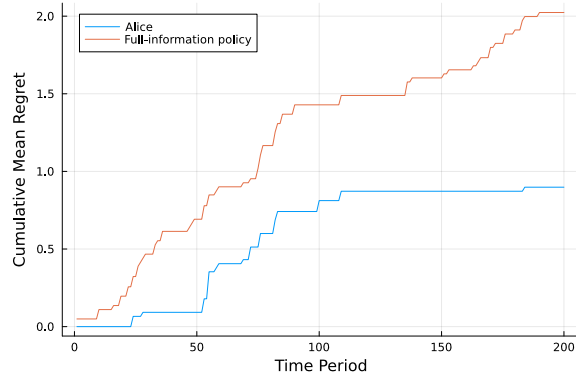


Figure 3: Cumulative mean regret under the learned policy (blue) compared to the full-information baseline (orange), for an example with τ having binary support.

4 Discussion

Our approach to algorithmic choice architecture uses model-based reinforcement learning to solve the information design problem of a sender deciding which options and features to show to a boundedly rational decision-maker. In order to model the effect of information overload and cognitive costs, we introduce a novel information-theoretic model of multi-attribute choice based on rational inattention theory. How robust the sender’s solution is to misspecified models of the receiver’s behavior is an open question.

By inferring preferences from boundedly rational behavior, our methodology can explicitly optimize for revealed consumer welfare, and is thus less paternalistic than traditional nudging. It also promises better alignment of artificial agents with the preferences of boundedly rational humans. This has implications beyond economics, for example for alignment research in artificial intelligence.

While the exposition has, motivated by benevolent nudging, focused on the case of aligned preferences, this need not be the case in our model: Alice might also want to persuade Bob, so our model shares features of Bayesian persuasion with a rationally inattentive receiver [cf. Bloedel and Segal, 2021], albeit with a more applied focus on situations in which receiver preferences are unknown and the sender is restricted to revealing true information.

Acknowledgments and Disclosure of Funding

Funding from the Max Planck Society and the Humboldt Foundation is gratefully acknowledged. The authors declare no competing interest.

References

- Alexander W. Bloedel and Ilya Segal. Persuading a Rationally Inattentive Agent. *Working Paper*, 2021.
- Stefan Bucher and Andrew Caplin. Inattention and Inequity in School Matching. *NBER Working Paper*, No. 29586, 2021.
- Stefan Bucher and Peter Dayan. Algorithmic Choice Architecture for Boundedly Rational Consumers. *In preparation*.
- Alison Bутtenheim, Robert Moffitt, and Alexandra Beatty, editors. *Behavioral Economics: Policy Impact and Future Directions*. The National Academies Press, Washington, DC, 2023. doi: 10.17226/26874.
- Andrew Caplin. The Science of Mistakes, Lecture Notes on Economic Data Engineering. *World Scientific Lecture Notes in Economics and Policy*, 2023. ISSN 2630-4872. doi: 10.1142/9789811262395_0005.
- Stefano DellaVigna and Elizabeth Linos. RCTs to Scale: Comprehensive Evidence From Two Nudge Units. *Econometrica*, 90(1):81–116, 2022. ISSN 0012-9682. doi: 10.3982/ecta18709.
- Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. The Challenge of Understanding What Users Want: Inconsistent Preferences and Engagement Optimization. *arXiv*, 2022. doi: 10.48550/arxiv.2202.11776.
- Maximilian Maier, František Bartoš, T D Stanley, David R Shanks, Adam J L Harris, and Eric-Jan Wagenmakers. No evidence for nudging after adjusting for publication bias. *Proceedings of the National Academy of Sciences*, 119(31):e2200300119, 2022. ISSN 0027-8424. doi: 10.1073/pnas.2200300119.
- Filip Matejka and Alisdair McKay. Rational Inattention to Discrete Choices: A New Foundation for the Multinomial Logit Model. *American Economic Review*, 105(1):272—298, 2015.
- Bartosz Maćkowiak, Filip Matějka, and Mirko Wiederholt. Rational Inattention: A Review. *Journal of Economic Literature*, 61(1):226–273, 2023. ISSN 0022-0515. doi: 10.1257/jel.20211524.
- Daniel McFadden. Choice - What Can Go Wrong? *NBER Working Paper*, No. 31165, 2023.
- Stephanie Mertens, Mario Herberz, Ulf J. J. Hahnel, and Tobias Brosch. The effectiveness of nudging: A meta-analysis of choice architecture interventions across behavioral domains. *Proceedings of the National Academy of Sciences*, 119(1):e2107346118, 2022. ISSN 0027-8424. doi: 10.1073/pnas.2107346118.
- David Silver and Joel Veness. Monte-Carlo planning in large POMDPs. *Advances in Neural Information Processing Systems*, 23, 2010.
- Christopher A. Sims. Implications of rational inattention. *Journal of Monetary Economics*, 50(3): 665–690, 2003. ISSN 0304-3932. doi: 10.1016/s0304-3932(03)00029-1.
- Michael Woodford. Modeling Imprecision in Perception, Valuation, and Choice. *Annual Review of Economics*, 12(1):579–601, 2020. ISSN 1941-1383. doi: 10.1146/annurev-economics-102819-040518.