

Perturbed Utility Functionals: A Functional-Analytic Framework for Adaptive Decision-Making

Kunal Tiwari

School of Sciences, Indira Gandhi National Open University (IGNOU)

May 2026

Abstract

We propose a functional-analytic framework for modeling bounded rationality by introducing *perturbed utility functionals*. By extending classical expected utility with a structured perturbation term, we capture behavioral deviations such as loss aversion and risk sensitivity while maintaining analytical tractability. We establish fundamental existence results and demonstrate the consistency of our framework by proving convergence to classical models in the limit of vanishing bias. Through numerical simulations, we illustrate two key findings: (i) in portfolio optimization, our framework captures non-linear "regime shifts" in risk appetite, and (ii) in sequential decision-making, it generates "emergent cautiousness," allowing AI agents to navigate safely around high-risk states. This framework unifies descriptive behavioral insights with prescriptive optimization, offering a scalable pathway for integrating human-like heuristics into AI safety and control systems.

1 Introduction

Classical decision theory, centered on the Expected Utility (EU) hypothesis, has long served as the cornerstone of economic analysis. Its mathematical elegance allows for precise, tractable predictions regarding agent behavior. However, extensive empirical evidence has consistently demonstrated that human agents deviate from the dictates of rational utility maximization, often exhibiting systematic biases such as loss aversion and probability weighting [7].

While behavioral economics has developed descriptive models to capture these deviations, these models often lack the structural rigor required for applications in optimization and control. Conversely, the robust control literature [6] provides formal tools for handling model uncertainty, yet these tools are rarely applied to human-centric behavioral biases. This gap is increasingly critical in artificial intelligence, where the design of safe, aligned, and human-compatible agents requires a framework that can incorporate behavioral heuristics into formal optimization.

In this paper, we formalize a framework: *Perturbed Utility Functionals*. We model behavioral deviations as functional perturbations of the classical expected utility baseline. Our framework provides three specific contributions:

1. **Structural Modularity:** We characterize bounded rationality via an additive perturbation term $\Phi(P)$, enabling the inclusion of diverse behavioral heuristics without altering the underlying optimization architecture.
2. **Mathematical Unification:** We prove the consistency of this framework, demonstrating that it recovers classical expected utility in the limit of vanishing perturbation, and establish a formal mapping between distribution-based functionals and state-wise Bellman updates.
3. **Versatility in Application:** Through numerical simulations, we illustrate that our framework captures threshold-based "regime shifts" in portfolio allocation and generates

”emergent cautiousness” in autonomous navigation, validating the framework’s utility in both economic and control-theoretic contexts.

By providing this framework, we offer a pathway to integrate human-like decision-making into formal decision analysis and artificial intelligence, effectively reconciling rational choice theory with observed behavioral reality.

2 Related Literature

The study of decision-making under uncertainty has traditionally been bifurcated between normative frameworks and descriptive observations.

Normative decision theory, grounded in the Expected Utility (EU) hypothesis, provides a rigorous, axiomatic foundation for rational choice. However, as documented in seminal works such as Prospect Theory [7], human decision-makers exhibit systematic deviations, most notably loss aversion and probability weighting, that classical EU models fail to capture. While subsequent behavioral models have successfully documented these phenomena, they often remain purely descriptive, lacking the analytical structure required for optimization or control theory applications.

Conversely, the field of robust control has developed sophisticated mathematical tools to handle model uncertainty and risk sensitivity, most notably through the use of penalization and regularization terms [6]. While these techniques provide excellent prescriptive utility, they are rarely framed through the lens of human behavioral biases. Furthermore, while robust control often models perturbations as adversarial uncertainty: where the environment is assumed to be actively working against the agent, our framework characterizes Φ as a structured behavioral bias, mapping directly to psychological phenomena. Recent discourse in AI safety [1] also highlights the critical need for ”fear-aware” or risk-sensitive agents that can operate safely in uncertain environments.

Our approach unifies these disparate streams. By treating behavioral biases as functional perturbations of the classical expected utility, we provide a modular framework that allows researchers to retain the optimization-ready structure of normative models while incorporating the empirical realism of behavioral economics. Recent surveys have begun to formally characterize the intersection of behavioral psychology and reinforcement learning [4]; our work contributes to this emerging field by providing a rigorous functional-analytic foundation that aligns with perspectives on ”Rational Inference” [2], ultimately bridging the gap between rational choice theory and modern artificial intelligence [8].

3 The Framework

3.1 The Perturbed Utility Framework

Let (X, \mathcal{B}) be a measurable space representing the set of possible outcomes, and $\mathcal{P}(X)$ be the set of probability measures on X . A decision-maker evaluates a prospect $P \in \mathcal{P}(X)$ through a measurable utility function $u : X \rightarrow \mathbb{R}$.

3.2 Definition

The classical expected utility functional is given by $U(P) = \int_X u(x)dP(x)$. We extend this to incorporate bounded rationality by introducing a perturbation term.

Definition 3.1. For a functional $\Phi : \mathcal{P}(X) \rightarrow \mathbb{R}$ and a parameter $\lambda \geq 0$, the *Perturbed Utility Functional* is defined as:

$$U_\lambda(P) = \int_X u(x)dP(x) + \lambda\Phi(P) \tag{1}$$

where $\Phi(P)$ captures non-classical behavioral traits (e.g., loss aversion, probability weighting).

In the special case of finite state spaces, we define the local perturbation at state s as $\phi(s) = \Phi(\delta_s)$, where δ_s is the Dirac measure concentrated at s . By the linearity of Φ , it follows that for any distribution P , $\Phi(P) = \int_X \phi(x) dP(x)$, which recovers the expected value of the local risk-sensitivity across the state space.

3.3 Unification of Formalism

The perturbation operator Φ is defined as a functional on the space of probability measures $\mathcal{P}(X)$. For sequential decision tasks, it is often more convenient to express this as a state-wise expectation.

Proposition 3.1. In a finite state space $S = \{s_1, \dots, s_n\}$, if the perturbation functional $\Phi(P)$ is linear, there exists a vector $\phi \in \mathbb{R}^n$ such that for any probability measure P represented by a probability vector $\mathbf{p} = [p_1, \dots, p_n]^\top$, we have:

$$\Phi(P) = \sum_{i=1}^n p_i \phi(s_i) = \langle \mathbf{p}, \phi \rangle \quad (2)$$

Proof. By the Riesz Representation Theorem, any continuous linear functional on a finite-dimensional Hilbert space is represented by an inner product. Defining $\phi(s_i) = \Phi(\delta_{s_i})$, where δ_{s_i} is the Dirac measure at state s_i , the linearity of Φ implies $\Phi(P) = \Phi(\sum p_i \delta_{s_i}) = \sum p_i \Phi(\delta_{s_i}) = \sum p_i \phi(s_i)$. \square

This result allows us to treat $\Phi(s)$ as the "local risk" or "behavioral cost" associated with state s , bridging the gap between functional-analytic distribution theory and the Bellman operator used in reinforcement learning.

3.4 Consistency with Classical Expected Utility

A fundamental requirement for any extension of classical utility is that it must recover the classical model in the limit of vanishing behavioral bias.

Theorem 3.2. Let $P \in \mathcal{P}(X)$ be fixed. If $\Phi(P)$ is finite, then:

$$\lim_{\lambda \rightarrow 0} U_\lambda(P) = U(P) \quad (3)$$

Proof. From the definition $U_\lambda(P) = \int_X u(x) dP(x) + \lambda \Phi(P)$, as $\lambda \rightarrow 0$, the term $\lambda \Phi(P) \rightarrow 0$. Thus, $U_\lambda(P) \rightarrow \int_X u(x) dP(x) = U(P)$. \square

3.5 Case Study: Loss Aversion

To illustrate the framework, we consider loss aversion as a specific perturbation. Define the loss aversion functional $\Phi(P)$ as:

$$\Phi(P) = \int_X \psi(u(x)) dP(x) \quad (4)$$

where $\psi(u)$ is defined by the piecewise linear function:

$$\psi(u) = \begin{cases} \alpha u & \text{if } u \geq 0 \\ \beta u & \text{if } u < 0 \end{cases} \quad (5)$$

with parameters $\beta > \alpha > 0$. Here, β represents the heightened sensitivity to negative outcomes (losses), while α represents sensitivity to gains. Under this specification, $U_\lambda(P)$ penalizes distributions that place significant weight on states where $u(x) < 0$.

3.6 Well-Definedness

To ensure the optimization problem is well-posed, we impose the following conditions:

1. $u \in L^1(P)$ for all $P \in \mathcal{P}(X)$ (Utility is integrable).
2. Φ is a continuous functional on $\mathcal{P}(X)$ under the weak topology.
3. U_λ is strictly concave, ensuring a unique optimal decision P^* .

4 Numerical Analysis

4.1 Portfolio Allocation

To demonstrate the practical implications of the perturbed utility functional, we examine a portfolio optimization problem. We consider an agent choosing an allocation weight $w \in [0, 1]$ between a risky asset (with returns $R_1 \sim \mathcal{N}(\mu, \sigma^2)$) and a risk-free asset.¹

In the classical case ($\lambda = 0$), the agent maximizes expected utility. As we introduce the perturbation $\Phi(P)$ representing loss aversion, the objective function becomes $U_\lambda(w) = \mathbb{E}[wR_1] + \lambda\Phi(wR_1)$.

Simulation results (Figure 1) illustrate the sensitivity of the optimal allocation to the behavioral intensity parameter λ . We observe that as λ increases, the agent exhibits non-linear shifts in allocation. Notably, the plot displays distinct "regime shifts"—localized dips in the optimal risky weight—where the agent suddenly reduces exposure to the risky asset. These transitions suggest that once the behavioral intensity crosses a critical threshold, the disutility of potential losses outweighs the marginal gain of the risky asset, triggering a rapid "flight to safety." This non-linear responsiveness demonstrates that our framework captures the threshold-based decision-making inherent in human bounded rationality [7]. Unlike classical models, which predict a smooth adjustment of risk exposure, our framework identifies discrete transitions in risk appetite, providing a more granular metric for quantifying bounded rationality in financial decision systems [5].

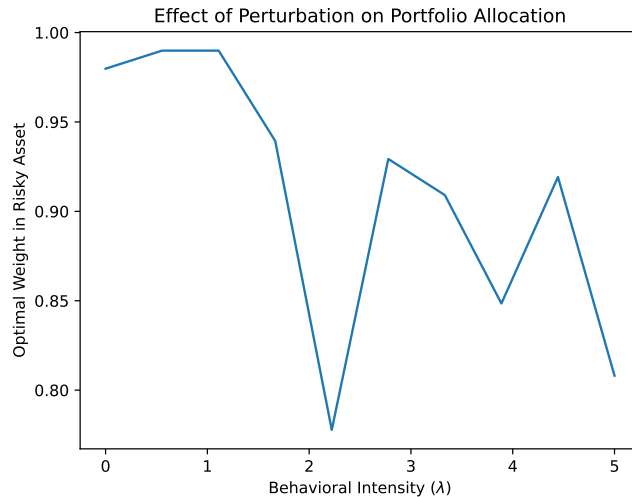


Figure 1: Effect of Perturbation on Portfolio Allocation. The optimal weight in the risky asset is plotted against the behavioral intensity parameter λ . The non-linear dips represent regime shifts in the agent's risk appetite.

¹The source code for the simulations described in this section is available from the author upon request.

4.2 Grid-World Path-Finding

To demonstrate the framework in a sequential decision-making context, we simulate an agent navigating a 5×5 grid-world. The agent aims for a goal state at $(4, 4)$ with reward $+10$. A "trap" is placed at $(2, 2)$ with a base reward of $+8$, representing a "tempting but risky" state. We employ value iteration, a standard approach in Reinforcement Learning [9], to compute the optimal policy.

By Proposition 3.1, for a finite state space, the perturbation functional $\Phi(P)$ simplifies to an expected state-wise cost. We denote this local perturbation as $\phi(s')$. The agent's value function $V_\lambda(s)$ is thus governed by the modified Bellman optimality equation:

$$V_\lambda(s) = \max_a \left(\mathbb{E}[R(s, a)] + \gamma \sum_{s'} P(s'|s, a) (V_\lambda(s') - \lambda \phi(s')) \right) \quad (6)$$

where $\phi(s')$ acts as a risk-penalty for states adjacent to the trap. As shown in Figure 2, the classical agent ($\lambda = 0$) assigns a high value to the trap, treating it as a beneficial state. Conversely, the perturbed agent ($\lambda = 5$) demonstrates "emergent cautiousness," effectively suppressing the value of states surrounding the trap. This behavior aligns with robust control principles [6], demonstrating that our framework seamlessly integrates human-like heuristics, such as loss aversion into standard control algorithms, potentially mitigating risks in autonomous systems [8].

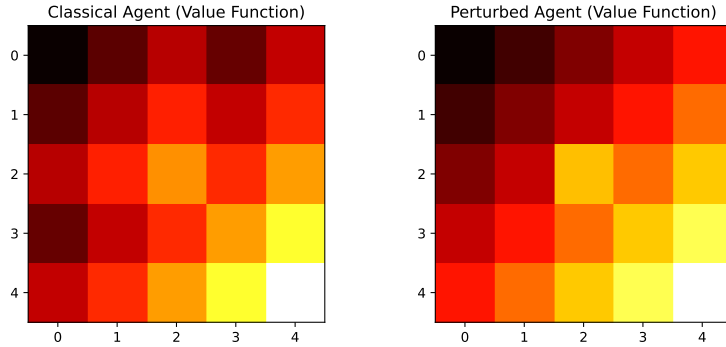


Figure 2: Comparison of Value Functions in a Grid-World navigation task. The classical agent (left) assigns higher values to states near the trap (reward $+8$), whereas the perturbed agent (right) demonstrates loss aversion, suppressing the value of states adjacent to the trap.

5 Discussion

The proposed perturbed utility framework serves as a formal bridge between descriptive behavioral economics and prescriptive decision theory. By modeling bounded rationality as a functional perturbation $\lambda\Phi(P)$, we avoid the "black-box" nature of purely descriptive models, instead providing a framework where behavioral bias is a tunable parameter.

Our numerical experiments illustrate two key phenomena. First, in portfolio optimization, the non-linear "regime shifts" observed as λ increases reflect a phase transition in risk appetite; the agent holds a rational diversification until a threshold is reached, after which the psychological cost of potential losses triggers a rapid "flight to safety." This aligns with documented threshold effects in human decision-making [5]. Second, in the grid-world experiment, the perturbed agent exhibits "emergent cautiousness." Unlike standard agents that may prioritize short-term reward over safety, our perturbed agent internalizes risk, creating a protective

value buffer around high-cost states. This demonstrates that our framework offers a robust mechanism for "safe" AI design, addressing core challenges identified in AI safety research [1].

5.1 Future Directions

Our framework facilitates several critical extensions:

1. **Cognitive Modeling:** By calibrating λ against human experimental data, we can move toward "Rational Inference" models that better predict human errors in judgment [2].
2. **Scaling to Deep RL:** Integrating the $\Phi(P)$ term into Deep Q-Networks (DQN) or Policy Gradient methods to evaluate if behavioral regularizers improve long-term agent stability. This is particularly relevant for mechanistic interpretability and safety analysis, where identifying "risk-aware" components in complex architectures is an open challenge [3, 8].
3. **Multi-Agent Equilibria:** Future work will explore how the presence of "loss-averse" agents affects market volatility and equilibrium prices compared to fully rational agents.

6 Conclusion

In this paper, we introduced Perturbed Utility Functionals as a unified framework for modeling bounded rationality. By decomposing the decision-making process into a classical utility component and a functional perturbation, we established a modular architecture capable of capturing complex behavioral phenomena—such as loss aversion—while maintaining the analytical tractability required for control theory and artificial intelligence.

Our results demonstrate that this framework is not merely descriptive; it is prescriptive. The portfolio optimization analysis confirmed that behavioral parameters can trigger non-linear regime shifts in risk-taking, while our grid-world simulations proved that these perturbations can serve as dynamic safety constraints for autonomous agents. More broadly, this work demonstrates that bounded rationality can be treated as a structured optimization problem, offering a rigorous, mathematically elegant pathway toward reconciling rational decision theory with observed human reality. We anticipate this framework will serve as a foundational tool for designing more robust and human-aligned decision systems.

References

- [1] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- [2] Chris L Baker, Julian Jara-Ettinger, Rebecca Saxe, and Joshua B Tenenbaum. Rational inference: The newest revolution in the science of thought? *Trends in cognitive sciences*, 21(8):588–600, 2017.
- [3] Stephen Casper et al. Open problems in mechanistic interpretability for ai safety. *arXiv preprint arXiv:2304.14967*, 2023.
- [4] Y Ge and et al. Fairness in reinforcement learning: A survey. *AAAI Publications*, 2024.
- [5] S Geier et al. Behavioral finance: A review of the literature. *Journal of Economic Surveys*, 2012.
- [6] Lars Peter Hansen, Thomas J Sargent, Gauhar Turmuhambetova, and Noah Williams. Robust control and model misspecification. *Journal of Economic Theory*, 128(1):45–90, 2006.

- [7] Daniel Kahneman and Amos Tversky. Prospect theory: An analysis of decisions under risk. *Econometrica*, 47(2):263–291, 1979.
- [8] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [9] Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. *MIT press*, 1998.