

Beyond Rational Choice: A Generative AI Framework for Modelling Cognitive Biases in Economic Decision Making

Gaurika Bhatia¹, Dr Nishant Kumar Singh²

¹Student, Department of Computer Science and Engineering, SRM Institute of Science and Technology

²Associate Professor, Department of Computer Science and Engineering, SRM Institute of Science and Technology

Abstract

When faced with uncertainty, people often diverge from traditional rational-choice theory due to bias-induced cognition and reasoning. Behavioural economics literature reveals that individuals frequently deviate from ideal decision-making methods which leads to systematic biases in their judgment, risk attitudes and preferences. While developments in Artificial Intelligence (AI), particularly in generative models, have led to better predictions of economic behavior, the majority of models continue to be based on rational assumptions and do not fully capture these biases. This paper proposes a generative AI approach to model biased decision-making processes. This approach integrates insights from behavioural economics and state-of-the-art generative models to capture a decision-making process that mirrors human behavior. It is designed to provide a range of decision pathways influenced by biases and compare them with the rational decision-making benchmark. This enables a more detailed understanding of the effects of different biases on decision-making in various scenarios. The mechanism captures major cognitive biases - such as anchoring, loss aversion, overconfidence, and confirmation bias - by modelling them as computational operations on structured decisions. This makes it capable of modelling non-linear and context-specific patterns. It also allows comparisons with measures of utility and intuitive visualisations of decision biases. The results demonstrate that models that incorporate cognitive biases are able to explain human behaviour that rational models cannot. This demonstrates the need for incorporating human thought processes in AI decision making. Overall, this work provides a scalable and intuitive framework for modeling human decisions, enabling construction of future AI decision systems that more closely align with expected human behaviour.

Keywords: Cognitive Bias-Aware Decision Making, Generative AI Modelling, Behavioral Economics, Human-Centric AI, Decision Modelling

INDEX TERMS

Artificial Intelligence, Generative AI, Behavioral Economics, Cognitive Bias Modeling, Decision-Making Systems, Explainable AI, Loss Aversion, Anchoring Bias, Overconfidence Bias, Confirmation Bias, Utility-Based Evaluation, Human-Centered AI.

INTRODUCTION

Economic decision-making has traditionally been seen as reliant to the notion that people make rational decisions by formally processing information to achieve satisfactory outcomes. This is the underpinning of classical and neoclassical economics, as well as impacting upon computational modelling of human decision-making. So we summarise that there is now enough empirical evidence in behavioural economics showing that we do not behave rationally in decision-making as we are biased. These include loss aversion bias, anchoring, over-confidence and confirmation which affect our view of the choices, risk and uncertainty. With advances in Artificial Intelligence (AI) technology and in particular in data-driven and generative modelling, there has been a growing tendency to build systems to model and predict human decision-making. So we can say that because of their ability to recognise patterns and optimise changes today's AI systems can play an important role in decision-making. Current research in this area still adheres to the ideals of rational decision-making. These systems emphasise on truth and expected utility but ignore the uncertain, dissimilar and irrational decision-making. This contributes to the mismatch between the predictions of AI systems and humans, and limits the use of these AI systems for real world applications, such as finance, consumer analysis and policy making usefulness of the systems in fields such as finance, consumer analytics and public policy

A. Research Problem and Question

The research problem investigated in this paper is the lack of currently available AI-driven decision models that adequately depict human decision-making in the presence of cognitive biases. Although there are recent developments in machine learning and generative decision models, but these models do not have ways to either: model how decisions change with different psychological states.

This presents the research question:

How can generative Artificial Intelligence (AI) be used to simulate and predict decisions influenced by cognitive biases in economic settings while being interpretable and aligned with behavioral theory?

To answer this question, it is necessary for us to understand how to generate decisions while also modelling the behavioural processes that lead to biases.

B. Limitations of Rational AI Models

The reliance on rational assumptions in AI-based systems creates a significant gap between computational predictions and actual human behavior. Traditional models process structured inputs such as probabilities and rewards but do not account for subjective perception or psychological influence. In real-world scenarios, individuals often overestimate gains, underestimate risks, or rely on heuristics shaped by prior experiences and beliefs.

Furthermore, most machine learning models are designed for predictive accuracy rather than behavioral simulation. While they can estimate likely outcomes based on historical data, they lack the ability to generate alternative decision pathways under varying cognitive conditions. This restricts their applicability in environments where understanding behavioral variability is critical when influenced by different cognitive biases

C. Challenges in Computational Modeling of Cognitive Bias

The computational modelling of cognitive biases posed some problems. Biases are complex, context-specific, and non-linear, making them a challenge to model using conventional mathematical techniques. Current approaches either reduce cognitive biases to fixed rules or neglect them altogether, leading to unrealistic and non-adaptive models.

Moreover, many AI models are "black-box" approaches, offering little insight into the decision-making

process. Explainability is crucial for building confidence, evaluating, and using decision-making models. Lack of interpretability also hampers applications of these models in decision-making environments.

D. Research Gap and Motivation

Despite increasing interest in combining behavioral insights with AI, there is a lack of unified frameworks that effectively integrate cognitive bias modeling with modern generative techniques. Existing approaches tend to fall into isolated categories, either focusing on theoretical behavioral analysis or purely data-driven prediction models. Very few studies attempt to simulate decision-making processes by explicitly incorporating multiple cognitive biases within a structured generative framework.

These limitations highlight the need for a framework that can model decision-making in a manner that is both computationally robust and behaviorally realistic.

Existing Work	Limitation
Behavioral Economic models	Descriptive, not computationally implementable
Traditional Machine Learning Model	Assume rational decision-making
Rule-based bias modeling	Limited scalability and flexibility
Black-box AI models	Lack interpretability and transparency

E. Proposed Contribution and Significance

This paper proposes a generative AI-based framework for modeling cognitive biases in economic decision-making. The framework integrates behavioral economic principles with generative modeling techniques to simulate both rational and bias-influenced decisions. By representing cognitive biases as computational operators, the system generates multiple decision pathways for a given scenario and enables structured comparison between optimal and behavior-driven outcomes.

Unlike conventional approaches, the proposed framework emphasizes both realism and interpretability. It allows users to trace how specific biases influence decision outputs, thereby providing deeper insights into the underlying cognitive processes. The modular design ensures scalability and adaptability across different application domains.

The significance of this research lies in its contribution toward the development of human-centered AI systems that more accurately reflect real-world behavior. By bridging the gap between behavioral theory and computational modeling, this work enables more reliable decision-support systems in domains where human behavior plays a critical role. It also establishes a foundation for future research in bias-aware AI, interpretable generative systems, and behaviorally informed decision modeling.

LITERATURE SURVEY

We can broadly categorize most previous research relevant to this work into four broad areas; behavioral economics, computational models of decision-making, explainable artificial intelligence and generative artificial intelligence. We will assess these areas here, and highlight their shortcomings, which justify our proposed approach.

A. Models of Decision-Making

Traditional economic theories rest on the framework of rational choice decision-making, in which indiv-

Individuals make decisions to optimise their expected outcomes. This approach was revised by Kahneman and Tversky, who incorporated the notion of cognitive biases and heuristics. Prospect Theory also revealed that people weigh outcomes in relation to a reference point and demonstrate loss aversion, with losses being more pronounced than gains. However, and while these models offer a strong theoretical account of human decision making, they are largely descriptive and not formulated in a computational manner. As a result, their use in contemporary AI is limited as they do not provide efficient ways to model decision processes under different psychologies.

B. Computational Models of Decision Making

Machine learning has made it possible to develop computer decision models that seek to learn from past outcomes. Such models use statistical learning, optimization methods, and identified patterns to model behavior in structured settings. But these models assume fixed and rational decision-making behaviours, which are not characteristic of human behaviour. Finally, the majority of computational models are not aimed at simulation. They do not spawn alternative decision paths based on different cognitive factors, and therefore fail to model variability and uncertainty in human decision-making.

C. Explainable Artificial Intelligence for Decision Systems

The need for Explainable Artificial Intelligence (XAI) is increasingly prominent in decision-making systems that require transparency and explanation. Studies in this field stress the importance of models which offer underlying reasoning and justifiable outcomes, especially in critical applications. However, current XAI approaches are generic and not specifically designed for behavioural models. They concentrate on prediction interpretation rather than psychological factors that affect decision-making. This results in a disconnect between explainable AI and cognitive bias modeling in AI decision-making systems.

D. Generative Models Based on Transformers

Generative AI has been revolutionized by transformer-based models due to their adaptive representation learning and powerful generation capabilities. This class of models has shown promising results for natural language comprehension and generation. But they can't be applied in decision modeling due to their black-box nature. Although they can generate meaningful patterns, they lack explicit representations of behavioural aspects like biases. This restricts their use to realistically simulate human behavior as the output lacks expressive power to explain controlled behavioral variations.

E. Identified Research

Gap Despite advancements in many fronts, some limitations remain: • Behavioral models offer theoretical guidance but are not computationally scalable • Machine learning models assume rationality and are predictive (not simulation) • Explainable AI does not support behavioral modeling • Generative models do not explicitly model biases These types of issues suggest a gap in current research: the lack of an integrated approach to combining behavioral realism, generative models and explainability. We require models that can achieve realistic decision making under different cognitive bias states, with explainable and interpretable results. The research clearly shows that no solution currently exists that combines realism of behavior, generative capacity and interpretability. This approach seeks to overcome this limitation by establishing a framework to combine modeling cognitive bias with generative AI in a controlled and interpretable manner. Further, recent studies suggest that contemporary AI models can unintentionally inherit and propagate potential biases from their training data, rather than simulating them in a safe and transparent way. This adds to the complexity of disentangling this simulated cognitive bias from potential algorithmic bias, making it more difficult to ensure safe and ethical use of these systems. Another key

factor is the lack of comparison models. Current systems simply provide the prescribed decision without comparison. Comparative analysis is limited in terms of a systematic comparison across rational decisions and biased decisions, which limits the prospect to delve into behavioral deviations and their effects. In addition, behavior theories are not connected with other levels of abstractions. Psychology is at the conceptual level, machine learning at the statistical level and generative AI at the representation level. But there is no common layer to bridge these levels into a system to model, generate and explain human decision-making in a structured manner.

1 Algorithm 1

Generative AI Framework for Cognitive Bias-Aware Economic Decision Modeling

Input Set of decision scenarios:

$$S = \{s_1, s_2, s_3, \dots, s_n\} \quad (1)$$

where each scenario represents an economic decision context defined by attributes such as risk level, expected return, uncertainty, and contextual constraints.

Output Rational decision outputs D_r , bias-influenced decision outputs D_b , and quantitative analysis of behavioral deviations.

Procedure

1.1 Scenario Ingestion Input decision scenarios s_i through an interactive interface (e.g., API or Streamlit-based system).

1.2 Feature Extraction and Structuring Extract structured attributes (risk, reward, uncertainty, temporal factors) to obtain raw representations S_{raw} .

1.3 Data Normalization and Alignment Clean and normalize input data to remove inconsistencies and produce standardized representations S_{clean} .

1.4 Contextual Feature Encoding Encode scenarios into high-dimensional embeddings E_i using modern representation techniques such as transformer-based encoders or tabular deep learning models (e.g., TabTransformer, FT-Transformer).

1.5 Cognitive Bias Injection Layer Apply parameterized cognitive bias operators B_k (loss aversion, anchoring, overconfidence, confirmation bias) to generate modified representations:

$$E_i' = B_k(E_i) \quad (2)$$

These operators simulate behavioral distortions through nonlinear transformations of feature importance and perception weights.

1.6 Generative Decision Modeling Define a generative function $G(\cdot)$, implemented using advanced architectures such as:

- Transformer-based generative models
- Diffusion-based decision models
- Reinforcement learning with human feedback (RLHF)-aligned generators

1.7 Rational Decision Generation Generate baseline rational decisions:

$$D_r = G(E_i) \quad (3)$$

1.8 Bias-Conditioned Decision Generation Generate bias-influenced decisions:

$$D_b = G(E_i') \quad (4)$$

1.9 Deviation Quantification Compute divergence between D_r and D_b using utility-aware metrics.

1.10 Multi-Bias Interaction Analysis Evaluate the combined effect of multiple biases and identify dominant behavioral drivers influencing decision outcomes.

1.11 Visualization and Interpretation Generate interpretable outputs including decision comparison graphs, bias impact heatmaps, and scenario-level behavioral summaries.

1.12 Explainability Layer Provide transparent reasoning by mapping decision deviations back to specific bias transformations and feature contributions using attribution techniques (e.g., SHAP, attention scores).

CONCLUSION AND FUTURE SCOPE

The growing integration of Artificial Intelligence into economic and decision-centric systems has significantly advanced predictive and analytical capabilities. However, most existing computational models continue to rely on rational decision-making assumptions, which fail to capture the inherent variability and psychological complexity of real-world human behavior. This limitation restricts the applicability of AI systems in domains where decisions are influenced by perception, uncertainty, and cognitive bias.

To address this challenge, this work introduced a generative AI-based framework for modeling cognitive biases in economic decision-making. The proposed approach combines principles from behavioral economics with modern generative modeling techniques to simulate both rational and bias-influenced decisions within a unified computational structure. By representing cognitive biases as parameterized transformations, the framework enables the generation of multiple decision trajectories for a given scenario and facilitates systematic comparison between normative and behaviorally realistic outcomes.

The experimental analysis demonstrates that incorporating biases such as loss aversion, anchoring, overconfidence, and confirmation bias produces decision patterns that more closely resemble human behavior. The results highlight that even minor variations in perception and weighting of information can lead to significantly different outcomes, emphasizing the importance of behavioral factors in decision modeling. The framework further enables structured analysis of these deviations, providing insights into how specific biases influence decision dynamics across different contexts.

A key contribution of this work lies in its emphasis on interpretability. Unlike conventional black-box systems, the proposed framework exposes intermediate representations, bias transformations, and decision comparisons, allowing users to trace how outputs are generated. This level of transparency enhances trust and supports practical adoption in applications such as financial decision support, consumer behavior analysis, and policy evaluation. Additionally, the modular architecture ensures flexibility, enabling the integration of new biases, alternative generative models, and domain-specific adaptations.

Despite these contributions, certain limitations must be acknowledged. The current framework primarily relies on simulated decision scenarios due to the scarcity of large-scale, high-quality datasets that explicitly capture cognitive bias in real-world settings. Furthermore, while the implemented bias operators provide meaningful approximations, they simplify complex psychological processes that are often context-dependent and nonlinear. The evaluation is also conducted in controlled environments, which may not fully reflect the diversity and unpredictability of real-world decision-making behavior.

Future research will focus on extending the framework in several directions. First, the integration of real-world behavioral datasets, including financial transaction data and experimental decision logs, can enhance the empirical validity of the model. Second, the incorporation of advanced generative architectures such as diffusion-based models, reinforcement learning with human feedback (RLHF), and domain-adaptive transformer systems can improve the diversity and robustness of generated decision outputs. Third, the development of adaptive and personalized bias models can enable the system to learn

user-specific behavioral patterns over time, leading to more accurate and context-aware decision simulations.

In addition, future work may explore multi-agent decision environments where interactions between individuals introduce collective behavioral dynamics, as well as the integration of explainable AI techniques that align model reasoning with human cognitive understanding. Expanding the framework toward real-time decision support systems and human-in-the-loop applications also represents a promising direction for practical deployment.

In conclusion, this research demonstrates that integrating generative AI with behavioral economic principles provides a viable pathway toward more realistic and human-centered decision modeling. By bridging the gap between rational computation and psychologically grounded behavior, the proposed framework contributes to the development of next-generation AI systems capable of understanding, simulating, and supporting complex human decision-making processes in dynamic environments.

REFERENCES

1. Kahneman D., Sibony O., Sunstein C.R., “Noise: A Flaw in Human Judgment”, Little Brown Spark, 2021.
2. Thaler R.H., “Behavioral Economics: Past, Present, and Future”, *American Economic Review*, 2020, 110 (12), 3569–3579.
3. Russell S.J., Norvig P., “Artificial Intelligence: A Modern Approach”, Pearson, 2021.
4. Brown T.B., et al., “Language Models Are Few-Shot Learners”, *Advances in Neural Information Processing Systems*, 2020.
5. Wei J., et al., “Chain-of-Thought Prompting Elicits Reasoning in Large Language Models”, *Advances in Neural Information Processing Systems*, 2022.
6. OpenAI, “GPT-4 Technical Report”, 2023.
7. Bengio Y., et al., “Deep Learning for Artificial Intelligence”, *Communications of the Association for Computing Machinery*, 2021, 64 (7), 58–65.
8. Touvron H., et al., “LLaMA: Open and Efficient Foundation Language Models”, 2023.
9. Vaswani A., et al., “Attention Is All You Need”, *Advances in Neural Information Processing Systems*, 2017.
10. Sutton R.S., Barto A.G., “Reinforcement Learning: An Introduction”, MIT Press, 2018.
11. Christiano P.F., et al., “Deep Reinforcement Learning from Human Preferences”, *Advances in Neural Information Processing Systems*, 2017.
12. Ouyang L., et al., “Training Language Models to Follow Instructions with Human Feedback”, *Advances in Neural Information Processing Systems*, 2022.
13. Ribeiro M.T., Singh S., Guestrin C., “Why Should I Trust You? Explaining the Predictions of Any Classifier”, *Knowledge Discovery and Data Mining Conference*, 2016.
14. Lundberg S.M., Lee S.I., “A Unified Approach to Interpreting Model Predictions”, *Advances in Neural Information Processing Systems*, 2017.
15. Lipton Z.C., “The Mythos of Model Interpretability”, *Queue*, 2018, 16 (3).
16. Doshi-Velez F., Kim B., “Towards a Rigorous Science of Interpretable Machine Learning”, 2017.
17. Guidotti R., et al., “A Survey of Methods for Explaining Black Box Models”, *ACM Computing Surveys*, 2019, 51 (5).

18. He K., et al., “Deep Residual Learning for Image Recognition”, Computer Vision and Pattern Recognition Conference, 2016.
19. Radford A., et al., “Improving Language Understanding by Generative Pre-Training”, OpenAI, 2018.
20. Ho J., et al., “Denoising Diffusion Probabilistic Models”, Advances in Neural Information Processing Systems, 2020.
21. Camerer C.F., “Behavioral Game Theory: Experiments in Strategic Interaction”, Princeton University Press, 2022.
22. Gigerenzer G., “Risk Literacy in Decision Making”, Annual Review of Psychology, 2023.
23. Wooldridge M., “A Brief History of Artificial Intelligence”, 2021.
24. Kahneman D., “Maps of Bounded Rationality”, American Economic Review, 2003.
25. Mullainathan S., Thaler R.H., “Behavioral Economics”, Handbook of Economics.
26. OpenAI, “GPT-4.5 and Multimodal Generative Systems: Advances in Reasoning and Alignment”, 2025.
27. Google DeepMind, “Gemini: A Family of Multimodal Large Models”, 2024.
28. Meta AI, “LLaMA 3: Open Foundation and Instruction-Tuned Models”, 2024.
29. Anthropic, “Claude Models: Constitutional AI and Alignment Techniques”, 2024.
30. Bai Y., et al., “Constitutional AI: Harmlessness from AI Feedback”, 2023.
31. Hoffmann J., et al., “Training Compute-Optimal Large Language Models”, DeepMind, 2022.
32. Perez E., et al., “Scaling Laws for Reward Model Overoptimization”, 2023.
33. Bubeck S., et al., “Sparks of Artificial General Intelligence: Early Experiments with GPT-4”, 2023.
34. Kojima T., et al., “Large Language Models Are Zero-Shot Reasoners”, Advances in Neural Information Processing Systems, 2022.
35. Bommasani R., et al., “On the Opportunities and Risks of Foundation Models”, Stanford Center for Research on Foundation Models, 2021.



Licensed under [Creative Commons Attribution-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-sa/4.0/)