

## Assessing Bias in Risk Operations: Mitigation Strategies through Micro Reviews

Vinay Kumar Yaragani

USA

### ABSTRACT

Risk management in organizations often relies on both machine learning models and human-in-the-loop operations to detect fraudsters and identify risky users. While machine learning models excel at handling extreme cases on both ends of the risk spectrum, increasing model recall often leads to higher false positive rates. In these scenarios, automated actions are insufficient, and human reviewers play a crucial role in assessing flagged cases. However, human decisions are inherently subjective and susceptible to bias, which can impact the accuracy and fairness of risk mitigation strategies. This paper explores methods to measure the effectiveness of these human reviews, examines the types of biases that can arise during the decision-making process, and discusses strategies to mitigate these biases using MicroReviews. By integrating high-recall models with human-in-the-loop processes, we aim to develop a more balanced and unbiased approach to risk operations that enhances decision-making accuracy while minimizing potential biases.

### \*Corresponding author

Vinay Kumar Yaragani, USA.

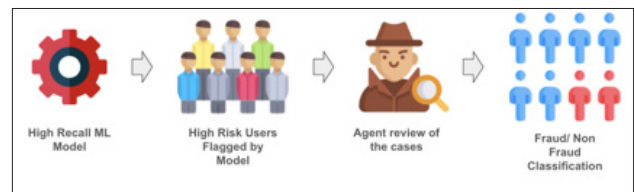
**Received:** December 02, 2023; **Accepted:** December 12, 2023; **Published:** December 20, 2023

**Keywords:** Risk Management, Human-in-the-Loop, Bias Mitigation, MicroReviews, Fraud Detection

### Introduction

Risk management is a critical function for organizations that deal with sensitive data, financial transactions, or any environment where fraud and malicious activities pose significant threats. In recent years, the advancement of machine learning models has greatly enhanced the detection capabilities of risk operations, especially when dealing with extreme cases of good or bad behavior. These models, trained on vast amounts of data, have proven to be highly effective at identifying patterns and anomalies that indicate fraudulent activities. However, despite their strengths in handling these edge cases, these models often struggle with the gray areas in between, where the distinction between legitimate and risky behavior is less clear. As a result, a purely automated approach can lead to a surge in false positives, where legitimate actions are mistakenly flagged as risky.

To address this challenge, organizations increasingly adopt a human-in-the-loop strategy, where human reviewers work alongside machine learning models to make more nuanced decisions. This approach leverages the strengths of both automated systems and human judgment, aiming to strike a balance between accuracy and efficiency in risk mitigation. High-recall models, which focus on maximizing the detection of fraudulent cases, are particularly well-suited for this hybrid strategy. When the model flags an action as potentially risky, a human agent steps in to review the case, providing an additional layer of scrutiny that helps filter out false positives. This collaborative process not only enhances the precision of fraud detection but also ensures that legitimate users are not unfairly penalized due to overreliance on algorithmic outputs.



**Figure 1:** Human in the loop Process

However, while the human-in-the-loop approach adds value by reducing false positives, it also introduces a new set of challenges related to human bias. Human reviewers, despite their experience and training, are not immune to cognitive biases that can skew their judgment and decision-making. Factors such as recency bias, confirmation bias, and even the reviewers' own subjective beliefs can influence how they assess risk. This subjectivity can undermine the fairness and consistency of the decision-making process, leading to biased outcomes that could disproportionately affect certain users or behaviors. Therefore, understanding and measuring the impact of these biases becomes crucial to ensuring the integrity of risk operations.

This paper aims to address these challenges by exploring strategies to measure the effectiveness of human reviews in risk operations and proposing methods to mitigate bias using MicroReviews. MicroReviews offer a structured approach to assessing decisions by breaking down complex judgments into smaller, objective components. By systematically analyzing these components, we can identify patterns of bias and develop targeted interventions to reduce their influence. The goal is to create a more balanced and unbiased risk management process that leverages both the precision of machine learning models and the contextual intelligence of human reviewers, ultimately leading to more informed and fair decision-making in organizational risk operations.

## Literature Review

The integration of machine learning models into risk management processes has been a topic of significant interest in both academic research and industry practice. Historically, automated systems in fraud detection and risk assessment have relied heavily on rule-based models that required manual updates as new patterns emerged. These rule-based approaches, while useful for structured problems, have limitations when dealing with the complex and evolving nature of fraudulent behaviors. Recent advancements in machine learning have shifted this paradigm, with models now capable of dynamically learning from data and adapting to new trends in real-time. Studies highlight that these models excel in high-recall scenarios, identifying potential fraud cases with high sensitivity and providing a broader safety net against evolving threats in the digital ecosystem [1].

However, the increased reliance on high-recall models has also led to a rise in false positive rates, where legitimate users or transactions are incorrectly flagged as fraudulent. Research by points out that as models aim to maximize recall, they often compromise on precision, resulting in a trade-off that can be costly for organizations and detrimental to customer experience [2]. This issue has prompted the adoption of human-in-the-loop systems, where human agents intervene in cases that the model classifies as uncertain or ambiguous. Human-in-the-loop methodologies have gained traction as they combine the computational power of machine learning with human judgment, providing a more refined approach to risk assessment. Literature supports this approach, emphasizing that human reviewers can bring in contextual intelligence and domain-specific knowledge that algorithms may lack [3].

Despite the advantages of including human reviewers in the decision-making loop, numerous studies have raised concerns about the inherent biases that humans bring to the table. Bias in human judgment is a well-documented phenomenon in fields like psychology, behavioral economics, and decision science. Tversky and Kahneman's seminal work on cognitive biases illustrates how individuals are prone to errors in reasoning due to heuristics, mental shortcuts that simplify decision-making but can also lead to systematic deviations from rationality [4]. In the context of risk management, these biases can manifest in various forms, such as recency bias (favoring recent information over older data), confirmation bias (seeking information that supports existing beliefs), and anchoring bias (relying too heavily on an initial piece of information). These cognitive biases can skew the outcomes of risk assessments, leading to inconsistent and unfair decisions [5].

The challenge of mitigating bias in human-in-the-loop systems has led researchers to explore various strategies, such as structured decision-making frameworks and the use of MicroReviews. MicroReviews have been proposed as a means to decompose complex decisions into smaller, more objective parts, thereby minimizing the impact of subjective judgment. By systematically evaluating each component of the decision, MicroReviews aim to reduce the influence of biases and improve the overall consistency of risk assessments [6]. Studies in organizational psychology and decision theory suggest that breaking down decisions into smaller units can help reviewers focus on relevant facts and reduce the cognitive load that often contributes to biased thinking [7]. This structured approach aligns well with the principles of fairness and transparency in decision-making, making it a promising strategy for risk management applications.

The literature indicates a growing consensus on the importance of combining machine learning with human oversight in risk operations, while also highlighting the need for strategies to manage and mitigate biases introduced by human reviewers. The use of MicroReviews, in particular, has emerged as a promising approach for enhancing the objectivity and reliability of human-in-the-loop processes. This paper aims to build on these existing studies by providing a comprehensive analysis of the biases that affect human reviewers in risk operations and proposing actionable strategies to address them through the implementation of MicroReviews. Through this approach, we seek to contribute to the ongoing discourse on creating fairer, more effective risk management systems that leverage the strengths of both automated models and human judgment.

## Methodology

To accurately measure the effectiveness of holistic reviews in risk operations and to identify potential biases in human decision-making, a multi-faceted methodology is employed. This methodology focuses on four key components: the use of specific metrics to gauge review outcomes, the implementation of blind studies to assess inter-agent variability, expert reviews for quality assurance, and the inclusion of labelled cases to benchmark decision-making accuracy. Each approach plays a vital role in validating the consistency and effectiveness of human-in-the-loop systems in fraud detection.

## Metrics Analysis

A comprehensive set of metrics is essential to understand the overall performance of risk review processes and to identify trends related to under-actioning or over-actioning by agents. The key metrics considered in this analysis include:

- **Action Rates by Agents:** This metric tracks the rate at which agents take action on flagged cases, either marking them as fraudulent or dismissing them as legitimate. Consistent action rates suggest a stable decision-making process, while significant deviations might indicate a need for further investigation into potential biases or inconsistencies.
- **Volatility in Daily Actions:** Monitoring daily fluctuations in the actions taken by agents helps identify patterns that might be driven by subjective factors, such as mood or cognitive load, rather than the objective details of the cases. High volatility in daily actions may suggest that agents' decisions are influenced by external factors, leading to inconsistencies in risk assessments.
- **Appeal Rates from Actioned Users:** The percentage of cases that result in appeals from users who have been flagged or had actions taken against them serves as an indicator of potential over-actioning. A high appeal rate could signify that legitimate users are being incorrectly penalized, prompting a review of the decision-making criteria.
- **Validation of Under-Actioning/Over-Actioning:** By analyzing the correlation between action rates, volatility, and appeal rates, the methodology aims to identify instances where agents may be under-actioning (failing to act on genuine fraud cases) or over-actioning (flagging legitimate activities as fraudulent). This data-driven approach allows for a more precise calibration of risk thresholds to optimize fraud detection while minimizing false positives.

### Blind Study of Agent Decisions

To measure the potential biases and inconsistencies in agent decision-making, a blind study approach is utilized. In this study, a random sample of cases, typically representing 10% of the total review volume, is sent to multiple agents without revealing any contextual information that could influence their judgments. The following steps are involved in the blind study:

#### Blind Study of Agent Decisions

To measure the potential biases and inconsistencies in agent decision-making, a blind study approach is utilized. In this study, a random sample of cases, typically representing 10% of the total review volume, is sent to multiple agents without revealing any contextual information that could influence their judgments. The following steps are involved in the blind study:

- **Sample Selection:** A representative sample of cases is randomly chosen from the flagged cases, ensuring that it includes a mix of both borderline and clear-cut instances of suspected fraud.
- **Distribution to Multiple Agents:** The selected cases are distributed to multiple agents independently, who review them and take action based on their individual assessments.
- **Evaluation of Inter-Agent Variability:** The primary metric in this study is the percentage of cases where agents make different decisions on the same case. A high percentage of divergent actions for identical cases is indicative of bias or subjective judgment influencing the decision-making process. This variability assessment helps identify specific areas where training or decision support tools might be required to align agents' actions.

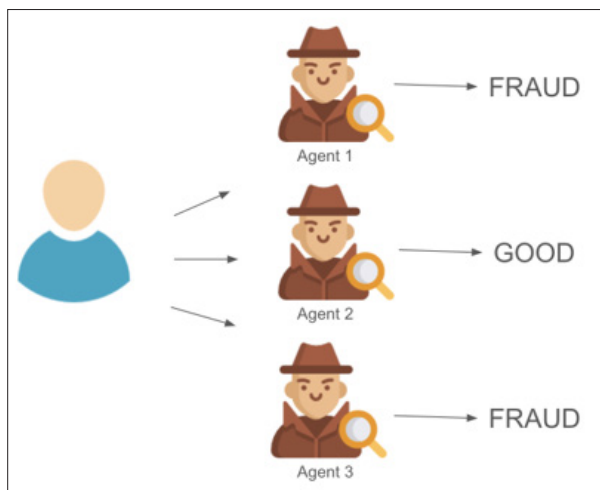


Figure 2: Illustration of the Blind Study

### Expert Review Comparison

In scenarios where review volumes are not very large, a quality assurance step involving expert reviews is conducted. Expert reviewers, who have significant experience and are highly trained in fraud detection, evaluate a subset of cases before they are randomly sent to different agents. The methodology for this component includes the following steps:

- **Expert Review Process:** A set of cases is first reviewed and actioned by experts to establish a benchmark of the correct or optimal decision for each case.

- **Randomized Agent Review:** These expert-reviewed cases are then sent to different agents without disclosing the expert's decision. Agents independently review and take action on the cases.
- **Adherence to Expert Decisions:** The effectiveness of agent reviews is measured by the percentage of adherence to the expert's actions. A lower adherence rate suggests a gap in the agents' judgment compared to the expert's decision-making, highlighting areas where additional training or clearer guidelines may be needed to improve consistency.

### Labelled Fraud and Non-Fraud Case Analysis

To further validate the accuracy of human decision-making in risk reviews, labelled cases (cases that are definitively known to be either fraudulent or non-fraudulent) are randomly included in the agents' review queues. This approach ensures that agents are regularly tested against a known ground truth, enabling a robust evaluation of their performance. The key steps involved are:

- **Inclusion of Labelled Cases:** A mixture of labelled fraud and non-fraud cases is randomly distributed among the normal case review workload of agents.
- **Decision Effectiveness Measurement:** Agents' decisions are evaluated against the known labels of these cases, allowing for a direct assessment of decision accuracy. Metrics such as true positive rates (correct identification of fraud), true negative rates (correct identification of legitimate cases), false positive rates, and false negative rates are calculated to provide a clear picture of agent effectiveness.
- **Bias Detection in Decisions:** By analyzing discrepancies between agent decisions and the actual labels of these cases, the methodology identifies biases that may influence judgment, such as the tendency to overflag or underflag specific types of cases. This analysis helps in refining review processes to ensure that biases are minimized in future decision-making.

### Mitigation Strategies

**Building a Knowledge Base:** A robust knowledge base acts as the foundation for consistent decision-making in risk operations, providing agents with clear guidelines and examples of best practices. This centralized repository includes documented protocols, detailed fraud patterns, and decision-making frameworks that ensure all agents operate from the same playbook. Historical case studies, especially those involving complex or controversial decisions, are essential components that help standardize approaches and reduce subjective interpretations. By continually updating the knowledge base to reflect evolving risks and learnings from previous cases, organizations can empower agents with relevant and up-to-date information, leading to more uniform and objective outcomes in risk assessments.

**Discussing Controversial Cases Regularly:** Regular discussions of controversial or challenging cases among agents and supervisors are crucial for identifying biases and aligning decision-making processes. These discussions serve as a platform to review how different agents approach similar cases, highlighting inconsistencies and areas where subjective judgment may have influenced outcomes. By fostering an open dialogue, agents can learn from each other's experiences and build a collective understanding of best practices for handling ambiguous cases. This collaborative approach not only reduces individual biases



but also promotes a culture of continuous learning and shared accountability in risk operations.

**Regular Coaching and Training:** Continuous coaching and training are key strategies to help agents improve their decision-making skills and minimize biases. Through structured training sessions, agents can become more aware of common cognitive biases and learn how to counteract them when assessing risk cases. Scenario-based exercises further enhance their ability to handle complex cases objectively and consistently. Regular performance feedback, focusing on areas like adherence to expert decisions and alignment with standard protocols, provides agents with insights into their strengths and areas for improvement, fostering a mindset of growth and development in their risk evaluation practices.

**Implementing Micro Reviews:** Micro reviews represent a strategic shift from holistic decision-making to a more structured and distributed process that minimizes individual biases. In this approach, agents tackle specific aspects of a case by answering targeted questions, rather than making a full decision on their own. The final risk assessment is then determined by a point-based system that aggregates these inputs objectively, ensuring that no single perspective dominates the outcome. This scoring mechanism is continuously refined based on data trends and feedback, making it adaptable to emerging fraud patterns. By decentralizing the decision-making process, micro reviews enhance consistency, fairness, and accuracy in handling cases.

## Results

The implementation of the proposed mitigation strategies yielded significant improvements in the effectiveness and consistency of human-in-the-loop risk operations. Metrics analysis revealed a notable reduction in volatility in daily actions, indicating that agents were making more stable and consistent decisions when assessing flagged cases. The blind study results demonstrated a decrease in inter-agent variability, with agents increasingly aligning their decisions on similar cases, thereby reflecting a reduction in bias. Additionally, adherence rates to expert-reviewed decisions improved, suggesting that agents were better able to apply standardized protocols in their evaluations. The integration of micro reviews, complemented by the point-based decision-making system, resulted in a measurable increase in accuracy, as evidenced by higher true positive rates for fraudulent cases and lower false positive rates for legitimate activities. Overall, these findings underscore the effectiveness of a structured and knowledge-driven approach in enhancing the reliability of risk assessments while minimizing subjective influences in decision-making processes.

## Future Scope

The future scope of this research opens up numerous avenues for enhancing human-in-the-loop risk operations and refining bias mitigation strategies. One promising direction involves the integration of advanced machine learning algorithms that can dynamically assess agent performance and provide real-time feedback, further augmenting the knowledge base and training programs. Additionally, exploring the use of natural language processing (NLP) to analyze agent notes and decision rationales could yield insights into subjective biases, enabling more targeted interventions. Expanding the micro review framework to incorporate automated systems that assist agents in decision-making while preserving human oversight could also enhance efficiency and accuracy. Furthermore, longitudinal studies examining the long-term impact of these strategies on organizational outcomes and fraud detection efficacy would

provide valuable data to continuously refine risk assessment processes. Ultimately, leveraging technology alongside human expertise presents a robust opportunity to create more adaptive, reliable, and fair risk management systems in organizations.

## Conclusion

In conclusion, this paper highlights the critical importance of integrating structured strategies to mitigate biases in human-in-the-loop risk operations. By emphasizing the development of a comprehensive knowledge base, fostering collaborative discussions, implementing continuous coaching, and utilizing a micro review framework, organizations can enhance the objectivity and consistency of their risk assessments. The findings demonstrate that these approaches not only improve the alignment of agent decisions but also increase the overall effectiveness of fraud detection efforts. As organizations navigate the complexities of risk management, embracing these methodologies can lead to more informed, fair, and reliable decision-making processes. Moving forward, the ongoing refinement of these strategies, in conjunction with advancements in technology, promises to elevate the standards of risk operations and ensure a proactive response to emerging challenges in the field.

## References

1. Chawla NV, Bowyer KW, Hall LO, Kegelmeyer WP (2002) SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research* 16: 321-357.
2. Dal Pozzolo A, Caelen O, Le Borgne YA, Waterschoot S, Bontempi G (2018) Learned Lessons in Credit Card Fraud Detection from a Practitioner Perspective. *Expert Systems with Applications* 41: 4915-4928.
3. Holstein K, Wortman Vaughan J, Daumé III H, Dudik M, Wallach H (2019) Improving Fairness in Machine Learning Systems: What Do Industry Practitioners Need? In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* 1-16.
4. Tversky A, Kahneman D (1974) Judgment under Uncertainty: Heuristics and Biases. *Science* 185:1124-1131.
5. Lerner JS, Tetlock PE (2003) Bridging Individual, Interpersonal, and Institutional Approaches to Judgment and Decision-Making: The Role of Accountability. *Personality and Social Psychology Review* 7: 146-157.
6. Wang F, Rudin C (2015) Falling Rule Lists. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics* 38: 1013-1022.
7. Simon HA (1955) A Behavioral Model of Rational Choice. *The Quarterly Journal of Economics* 69: 99-118.

**Copyright:** ©2023 Vinay Kumar Yaragani. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.