# AI Governance in Insurance: Establishing Controls for Transparent and Auditable Decisions

## Jalees Ahmad

jaleesahmad07@gmail.com

**Abstract:**

The insurance industry is undergoing a structural metamorphosis driven by the integration of artificial intelligence across the entire value chain. From automated underwriting and real-time risk assessment to sophisticated fraud detection and claims triage, AI-driven systems are replacing static, rule-based legacy processes. While these advancements offer significant improvements in operational efficiency and pricing precision, they also introduce systemic risks related to algorithmic bias, opacity, and regulatory non-compliance. This report examines the foundational requirements for establishing robust AI governance frameworks designed to ensure transparency and auditability. By analyzing global regulatory developments—including the European Union AI Act and the National Association of Insurance Commissioners (NAIC) Model Bulletin—this study delineates the technical and organizational controls necessary for responsible AI deployment. Key focus areas include the implementation of Explainable AI (XAI) techniques, the establishment of comprehensive "AIS Programs," the adoption of semantic record-keeping, and the integration of rigorous bias mitigation strategies. The findings emphasize that governance must move beyond a check-the-box compliance exercise toward a "defensible-by-documentation" architecture that preserves consumer trust while fostering sustainable innovation.

## INTRODUCTION

The traditional foundations of the insurance industry—grounded in human judgment, linear actuarial models, and historical risk pools—are being fundamentally reconfigured by the rapid ascent of artificial intelligence (AI) and machine learning (ML). This digital transformation represents a shift from a reactive "detect and repair" framework to a proactive "predict and prevent" model, wherein insurers leverage vast repositories of unstructured data to anticipate losses and personalize policyholder interactions. As insurers sit upon a "treasure-trove" of big data, the ability to utilize this resource through AI has become a primary driver of competitive advantage, enabling faster underwriting, instant claims processing, and more accurate risk segmentation.

However, the rapid diffusion of these technologies has created an accountability gap. The inherent complexity of modern AI architectures, particularly deep neural networks and ensemble methods, often results in "black-box" models where the logic underpinning a specific decision is obscured even from the developers themselves. In a highly regulated sector such as insurance, where decisions regarding premium rates or coverage denials carry significant socio-economic consequences, this lack of transparency is untenable. Regulators and consumer advocates are increasingly concerned that opaque algorithms may inadvertently codify societal biases, leading to unfair discrimination against protected classes without the possibility of human recourse or meaningful explanation.

The urgency for establishing robust governance controls is further underscored by the evolving global legal landscape. The European Union's AI Act and the NAIC Model Bulletin in the United States have established new precedents for algorithmic accountability, designating certain insurance applications as "high-risk" and demanding stringent documentation, data quality standards, and human oversight. For insurers, governance is no longer a peripheral technical concern but a core strategic imperative that influences everything from

capital efficiency under Solvency II to long-term reputational stability. This report provides a comprehensive analysis of the mechanisms required to establish transparent and auditable AI systems, ensuring that the "AI revolution" in insurance remains aligned with the principles of fairness, accountability, and legal compliance.

## Theoretical Foundations of AI Adoption in Insurance

The transition toward AI-centric insurance operations is not merely a technological upgrade but a shift that aligns with established theoretical models of organizational change and innovation. Understanding these foundations is critical for executives who must balance the pressure to innovate with the need for stable governance.

The willingness of insurance organizations to implement AI solutions is deeply influenced by perceptual factors described in the Technology Acceptance Model (TAM). When insurance leaders evaluate AI tools for sensitive functions like underwriting or claims processing, their perceptions of the technology's "usefulness" and "ease of implementation" often determine the speed and depth of adoption. Recent academic research has extended these models to account for the unique challenges of risk-sensitive environments, noting that the adoption of AI across the insurance landscape follows patterns consistent with the Diffusion of Innovation Theory. Technology-oriented firms typically pioneer these implementations, demonstrating measurable returns on investment (ROI) before practices spread to more traditional, risk-averse segments of the industry. This diffusion is particularly evident in the uneven adoption of machine learning algorithms across different market segments. While property and casualty (P&C) lines have aggressively adopted automated risk assessment for emerging threats like cyber-risk, life and health segments have been more cautious due to higher regulatory hurdles and the sensitive nature of the data involved. The current research landscape reveals that while AI yields substantial cost savings and efficiency gains, these benefits are most pronounced in "business-as-usual" regimes. Once data drift or novel threats emerge, the limitations of current AI capabilities—specifically the inability to reason on a global basis or understand causality—necessitate a continued reliance on human expertise.

## Mapping AI Applications and Their Governance Implications

The application of AI in insurance is pervasive, touching every stage of the lifecycle. However, each application carries a different risk profile, necessitating a proportionate approach to governance.

## Algorithmic Underwriting and Pricing

Underwriting 2.0 represents the next evolutionary phase in risk assessment, moving away from static, rule-based systems toward self-learning models. By analyzing complex datasets, including historical claims, real-time telematics, and unstructured text from medical reports, these systems can identify patterns that elude traditional actuarial methods. For instance, deep learning algorithms can process images of properties or vehicle damage to estimate risk more accurately than a manual survey.

The governance challenge here lies in "proxy discrimination." AI models must be designed to avoid using seemingly neutral variables—such as location or shopping habits—that may correlate highly with protected characteristics like race or religion. Research by Samiuddin et al. (2023) demonstrated the efficacy of Deep Neural Networks (DNN) in predicting medical insurance costs, outperforming traditional machine learning algorithms in both speed and accuracy. Yet, as Kaushik et al. (2022) noted, while these models can achieve accuracy rates of over 92%, the potential for "unjustified premium differentials" among protected groups remains a significant regulatory concern.

## Claims Management and Fraud Detection

In claims processing, AI acts as an accelerator for First Notice of Loss (FNOL) intake and triage. Agentic AI systems can intake multi-channel data, extract relevant details using NLP, and trigger workflows without human intervention. For straightforward claims, insurers like Lemonade have demonstrated "straight-through processing" where a genuine claim is validated and settled in as little as two seconds.

From a governance perspective, the risk shifts toward the "transparency of denial." If an automated system denies a claim, the insurer must be able to provide a clear, auditable rationale for that decision. Furthermore,

in the realm of fraud detection, AI uses predictive analytics to identify "red flags" and connect dots across disparate claims that might indicate a coordinated fraud ring. While this protects the risk pool, the governance framework must ensure that these "flags" do not unfairly target vulnerable populations due to biased training data.

| AI Application Area | Governance Risk Profile | Primary Control Mechanisms |
|---|---|---|
| **Underwriting & Pricing** | High: Direct impact on consumer access and cost. | Bias testing, fairness metrics, XAI (SHAP/LIME). |
| **Claims Triage** | High: Potential for wrongful denial of benefits. | Human-in-the-loop, appeals process, audit logs. |
| **Fraud Detection** | Medium/High: Risk of false positives and profiling. | Data lineage, pattern validation, manual investigation. |
| **Customer Support** | Low/Medium: Risk of misinformation or privacy breaches. | Transparency (AI disclosure), data minimization. |

**The Global Regulatory Landscape: Prescriptive and Principle-Based Models**
Governance frameworks are currently bifurcated between the prescriptive mandates of the European Union and the more principle-based guidance emerging in the United States and through international bodies like the OECD.

**The EU AI Act and Solvency II Integration**
The EU AI Act (Regulation (EU) 2024/1689) explicitly classifies AI systems used for risk assessment and pricing in life and health insurance as "high-risk". This classification triggers a comprehensive set of requirements, including the establishment of a quality management system, the maintenance of detailed technical documentation, and the implementation of robust data governance. High-risk systems must be designed for "record-keeping" to enable automatic event logging, and providers must ensure that the AI allows for effective human oversight.

A significant development in recent research is the linkage between AI fairness and regulatory capital under Solvency II. Academic studies have derived closed-form mappings from legal fairness metrics—such as statistical parity difference and disparate impact ratio—to changes in the Solvency II Solvency Capital Requirement (SCR). This implies that biased models are not just a legal liability but a financial one, as they can lead to increased capital charges and potential fines of up to 7% of worldwide annual turnover.

**The US NAIC Model Bulletin and State-Level Actions**
In contrast to the EU's legislative approach, the NAIC Model Bulletin on the Use of Artificial Intelligence Systems by Insurers provides a framework for state-level regulation. The bulletin is "prescriptive in nature" regarding the need for an "AIS Program," which must be a documented, written program for the responsible use of AI. Insurers are expected to vest responsibility for AI oversight with senior management accountable to the board of directors.
State-level adoption of this bulletin has been rapid. As of mid-2025, 24 states have adopted the model bulletin, while others like New York, Colorado, and California have enacted even more specific regulations. For example:
● **New York (Circular Letter No. 7)**: Demands robust oversight of AI used in pricing and underwriting, requiring insurers to demonstrate that their tools do not proxy for protected classes.

- **Colorado (SB21-169)**: Prohibits the use of external consumer data and predictive models in ways that result in unfair discrimination based on protected characteristics.
- **California (SB 1120)**: Specifically restricts health insurers from denying or modifying coverage based *solely* on an algorithm; any adverse determination must be reviewed by a licensed clinician.

## Establishing Structural Governance: The AIS Program

The central pillar of AI governance in the insurance sector is the "AIS Program," a comprehensive framework that embeds accountability and risk management into the organizational fabric.

## Accountability and Leadership Structure

The NAIC and EIOPA both emphasize that AI governance is not a mere technical concern but a leadership priority. An effective AIS Program must define clear lines of responsibility. Senior management is responsible for setting the insurer's AI strategy and ensuring that the organization possesses the "AI literacy" required to manage these systems effectively. The governance structure should ideally include stakeholders from actuarial, data science, underwriting, compliance, and legal departments.

## Model Risk Management (MRM) Adaptation

The frameworks introduced to deal with quantitative model risk in the 2010s are now being adapted for AI. However, traditional MRM is often insufficient because AI models are more dynamic and "digest mountains of data" differently than linear models. High-performing MRM programs must now include AI-specific safeguards:

- **Centralized Model Inventory**: Every AI system affecting pricing, underwriting, or claims must be cataloged and ranked by its risk exposure.
- **Rigorous Validation and Testing**: This includes performance testing to ensure the model behaves as designed and stress testing to reveal how it handles adverse scenarios beyond its stated scope.
- **Ongoing Monitoring**: Insurers must implement systems to test AI outputs for "model drift"—where the performance of a model degrades over time due to changes in real-world data patterns.

| AIS Program Component | Key Deliverables | Regulatory Expectation |
|---|---|---|
| **Governance Structure** | Board-approved AI Policy; Cross-functional AI Committee. | Clear accountability and strategic alignment. |
| **Risk Management** | Model Inventory; Risk Assessment Framework. | Identification and mitigation of foreseeable harms. |
| **Internal Controls** | Data Integrity Checks; Bias Audits. | Prevention of unfair discrimination and errors. |
| **Audit Function** | Internal and Third-party Audit Reports. | Independent verification of compliance and accuracy. |

## Data Governance and the Mitigation of Algorithmic Bias

Data is the "main ingredient" of AI success, but it is also the primary source of risk. Robust data governance is paramount to minimize the risk of biased or flawed AI outputs.

## Ensuring Data Quality and Representative Sets

Insurers must ensure that their training and testing datasets are relevant, accurate, and representative of the populations they serve. This is particularly challenging in "data-scarce" contexts or when using historical datasets that may reflect past societal prejudices. The EU AI Act requires that high-risk systems be trained on

data that is "free of errors and complete" to the best extent possible. In practice, this means insurers must validate source systems, resolve missing fields, and document data lineage before models are deployed.

**Quantifying Fairness and Addressing Bias**
Bias detection is not a monolithic task; it spans the entire lifecycle of model development. Insurers must integrate fairness metrics into the core fabric of AI design. One common metric is the **Disparate Impact Ratio (DIR)**, which compares the favorable outcome rates between a protected group and a reference group. The formula for the Disparate Impact Ratio can be expressed as:

$$DIR = \frac{P(\text{Outcome} = \text{Positive} | \text{Group} = \text{Unprivileged})}{P(\text{Outcome} = \text{Positive} | \text{Group} = \text{Privileged})}$$

Regulatory standards, such as the "four-fifths rule," often suggest that a ratio below 0.80 indicates potential discrimination. To mitigate identified biases, insurers are exploring strategies such as:
- **Fairness-aware data curation**: Modifying the dataset during the preprocessing stage to remove discriminatory patterns.
- **Adversarial Debiasing**: Training the AI model in a way that penalizes the system for being able to predict a protected attribute from the features.
- **Post-processing remediation**: Adjusting model outputs after they are generated to ensure they meet fairness criteria.

**Explainable AI (XAI) and the Quest for Transparency**
Transparency is essential for maintaining trust between insurers, customers, and regulators. Explainability is the mechanism by which the "black box" of AI is opened, allowing stakeholders to understand why a specific decision was made.

**Technical Approaches to Interpretability**
Insurers are increasingly moving toward "Algorithmic Underwriting 2.0," which incorporates XAI tools to provide auditable rationales for risk assessments. Techniques such as **LIME** (Local Interpretable Model-agnostic Explanations) and **SHAP** (Shapley Additive Explanations) have become industry standards.
- **SHAP** analysis allows underwriters to see the specific weight of each variable (e.g., body mass index, smoking status, or location) on a premium price.
- **LIME** provides a "local" explanation for a single individual's outcome, making it easier to explain a denial to a specific customer.

**Tailoring Explanations for Stakeholders**
Transparency is not a one-size-fits-all requirement. Governance frameworks must ensure that explanations are adapted to the recipient.
- **Authorities and Auditors**: Require global and comprehensive technical explanations of the system logic, including validation reports and error rates.
- **Consumers**: Must be informed when they are subject to an automated decision and provided with clarifications in "simple, clear, and non-technical language" regarding factors with a material impact on their outcomes.
- **Internal Users (Underwriters/Claims Adjusters)**: Need actionable insights that allow them to validate or override an AI recommendation with expert judgment.

**Auditability and the "Defensible by Documentation" Strategy**
Rigorous AI auditing is a strategic imperative to ensure sustained stakeholder trust. For insurers, this means building an AI system program that is "defensible by documentation," where every automated decision can be traced and tested.

**Semantic Record-Keeping and Data Lineage**
To satisfy regulators, insurers must capture more than just inputs and outputs. They must maintain a "semantic record" of the decision-making process. This includes:

- **Capturing Intermediate Logic**: Recording the retrievals and reasoning processes used by the AI before it reaches a final decision.
- **Data Lineage Tracking**: Building a branched data tree that shows exactly how data flowed from all source systems (e.g., broker emails or customer demographics) into the AI.
- **Version Pinning**: Ensuring that the organization can pin specific API versions and trace a decision back to the exact version of the model and prompt running at that time.

## Scenario-Based Auditing and Monitoring

Traditional model validation is often insufficient to address the dynamic nature of AI risks. Insurers are increasingly adopting "scenario-based auditing," which simulates real-world and edge-case situations to uncover vulnerabilities.

- **Bias Audits**: Systematic evaluation of algorithms for disparate impacts stratified by protected characteristics such as age, gender, and ethnicity.
- **Drift Monitoring**: Continuous evaluation of model performance across different time periods and populations to detect when a model's accuracy begins to degrade.
- **Audit Logs**: Tamper-resistant logs that record all system modifications, training data updates, and model validation outcomes.

| Audit Dimension | Requirement Detail | Objective |
|---|---|---|
| **Data Provenance** | Mapping sources and variable descriptions. | Verify data quality and integrity. |
| **Feature Selection** | Reviewing inputs for proxy variables. | Prevent discriminatory outcomes. |
| **Logic Traceability** | Semantic records of AI reasoning steps. | Enable independent review and appeals. |
| **Performance Logs** | Historical record of versioning and changes. | Ensure accountability over the lifecycle. |

## The Human-in-the-Loop: Accountability and Hybrid Models

Despite the trend toward automation, human judgment remains the lifeblood of the insurance industry. The shift is not toward replacing humans, but toward a "hybrid" model where AI serves as a "sparring partner" for human experts.

## Accountability and Final Judgment

A fundamental tenet of AI governance is that an algorithm cannot be legally or ethically accountable; only a human can. In complex underwriting scenarios or high-value claims, human underwriters and adjusters are essential to make final judgments with myriad factors in mind that an AI might miss. Senior leadership—the Administrative, Management, or Supervisory Body (AMSB)—holds the ultimate responsibility for AI use and must possess sufficient knowledge of the potential risks.

## Enhancing Professional Discretion

AI-driven tools can improve human performance by reducing the time spent on "grunt work," allowing senior professionals to focus on the most nuanced or exceptional cases. For example, Generative AI can condense lengthy medical reports into brief summaries, highlighting key risk areas so underwriters can focus on decision-making. This "human-on-the-loop" approach ensures that expert judgment can validate or override AI-generated decisions, providing a necessary safety net against algorithmic errors.

**CONCLUSION**

Establishing controls for transparent and auditable AI decisions is no longer a choice for the insurance industry but a condition of its survival in a digital-first regulatory environment. The transition to AI-driven operations offers immense promise for efficiency and precision, yet these benefits cannot be realized without a robust governance framework that addresses the unique risks of algorithmic opacity and bias.

This report has outlined the essential components of such a framework: a leadership-driven AIS Program, rigorous data governance centered on representative datasets and bias mitigation, and the technical implementation of Explainable AI and semantic record-keeping. The "defensible-by-documentation" strategy ensures that insurers can not only achieve regulatory compliance but also build a trust dividend with policyholders and partners. As AI systems assume increasingly influential roles in critical functions like life and health underwriting, the integration of human oversight and ethical principles into the AI lifecycle will be the defining factor in distinguishing the winners of the insurance digital transformation.

**REFERENCES:**

1. **Abraham, R., Schneider, J., & vom Brocke, J. (2023).** "Data Governance: A Conceptual Framework, Structured Review, and Research Agenda," *International Journal of Information Management*.
2. **Adeoye, O., et al. (2024).** "Artificial Intelligence Applications in Health Insurance: A Systematic Review," *Journal of Risk and Insurance*.
3. **Bhattacharya, S., Castignani, G., Masello, L., & Sheehan, B. (2025).** "AI revolution in insurance: bridging research and reality," *Frontiers in Artificial Intelligence*.
4. **Charpentier, A. (2024).** "Insurance and Big Data: The Ethical and Regulatory Challenges of Algorithmic Pricing," *Journal of Business Ethics*.
5. **Davis, F. D. (1989).** "Perceived Usefulness, Perceived Ease of Use, and User Acceptance of Information Technology," *MIS Quarterly*.
6. **Eling, M., & Lehmann, M. (2018).** "The Impact of Digitalization on the Insurance Value Chain and the Insurability of Risks," *The Geneva Papers on Risk and Insurance - Issues and Practice*.
7. **Gabelaia, S., et al. (2024).** "Technology Acceptance and the Deployment of Complex AI Systems in Risk-Sensitive Environments," *International Journal of Digital Transformation*.
8. **Kaushik, A., et al. (2022).** "Health Insurance Premium Prediction Using Machine Learning Regression Framework," *Academic Journal of Intelligent Systems and Machine Learning*.
9. **Kujala, J. (2023).** "Transparency and explainability of AI systems: From ethical guidelines to requirements," *Information and Software Technology*.
10. **Ma, Y., & Sun, L. (2020).** "Predicting Mortality Risk with Machine Learning: A Comparison of Traditional Actuarial Models and AI Algorithms," *North American Actuarial Journal*.
11. **Mahajan, S., Agarwal, R., & Gupta, M. (2025).** "Algorithmic Bias Under the EU AI Act: Compliance Risk and Capital Implications for Insurers," *Risks*.
12. **Orji, C., & Ukwandu, E. (2024).** "Explainable AI in Medical Insurance Cost Prediction: A Comparative Study of SHAP and ICE Plots," *World Journal of Advanced Research and Reviews*.
13. **Owens, E., Sheehan, M., et al. (2022).** "Explainable Artificial Intelligence (XAI) in Insurance," *Risks*.
14. **Puschmann, T. (2017).** "Fintech and the Insurance Industry: Exploring the Transformation of the Value Chain," *Business & Information Systems Engineering*.
15. **Rogers, E. M., et al. (2014).** *Diffusion of Innovations*, Free Press.
16. **Samiuddin, M., et al. (2023).** "Deep Learning for Health Insurance Cost Prediction: A Deep Neural Network Approach," *Central European Management Journal*.
17. **Shrestha, Y. R., et al. (2021).** "Algorithms and Decision Making in the Insurance Industry: Future Prospects and Challenges," *MIT Sloan Management Review*.
18. **Tewari, S. (2025).** "AI Powered Data Governance - Ensuring Data Quality and Compliance in the Era of Big Data," *Journal of Artificial Intelligence and Governance Systems*.