# RISK MANAGEMENT MAGAZINE

## Vol. 20, Issue 2
## May – August 2025

In collaboration with

## IN THIS ISSUE

### PAPERS SUBMITTED TO DOUBLE-BLIND PEER REVIEW

### NON REFERRED PAPERS

Scientific journal recognized by ANVUR and AIDEA

anvur

ACCADEMIA ITALIANA DI ECONOMIA AZIENDALE

**The authors bear sole responsibility for the opinions expressed in the articles.**

MAILED TO AIFIRM SUBSCRIBERS WHO ARE RESIDENT IN ITALY AND DULY REGISTERED

Journal printed on 27th August 2025

# RISK MANAGEMENT MAGAZINE

## Peer review process on papers presented for publication

The papers that are presented to our magazine for publication are submitted anonymously to a double level of peer review.

The first level is a review of eligibility, implemented on the paper by the members of the Editorial Board, who assess the adequacy of the paper to the typical topics of the magazine.

The second level is a review of suitability for publication, implemented on the paper by two referees, selected within the Editorial Board, the Scientific Committee or externally among academics, scholars, experts on the subject who assess the content and the form.

Email for article submission: risk.management.magazine@aifirm.it;

## Editorial regulation

"Risk Management Magazine" is the AIFIRM (Italian Association of Financial Industry Risk Managers) journal, fully dedicated to risk management topics.

The organization includes the managing editor, a joint manager and an Editorial Board and a Scientific Committee composed by academics.

The magazine promotes the diffusion of all content related to risk management topics, from regulatory aspects, to organizational and technical issues and all articles will be examined with interest through the Scientific Council.

The papers shall be presented in Microsoft Word format, font Times New Roman 10 and shall have between 5.000 and 12.000 words; tables and graphs are welcome.

The bibliography shall be written in APA format and shall accuratel specify the sources.

An Abstract in English is required (less than 200 words) highlighting the Key words.

The authors bear sole responsibility for the opinions expressed in the articles.

The Statement on ethics and on unfair procedures in scientific publications can be found on our website www.aifirm.it.

# Risk Culture & Culture Risk: not a play on words.

**Rosa Cocozza (Università degli Studi di Napoli Federico II) – Fernando Metelli (Albaleasing)**

**Corresponding Author: Rosa Cocozza (rosa.cocozza@unina.it)**

## Abstract

The purpose of this article is to suggest a primer for culture risk, aimed at outlining actionable and practical approaches distinguishing between 'risk culture' and 'culture risk'. The topic, originally addressed by the Financial Stability Board (FSB, 2014), has recently garnered renewed interest due to the Draft Guide on Governance and Risk Culture disseminated by the European Central Bank (ECB, 2024), setting out supervisory expectations, informed by the Capital Requirements Directive (CRD), European Banking Authority (EBA) guidelines, and international standards. Although the subject may be perceived as abstract, nevertheless it holds significant concrete relevance, despite the inherent challenges of measuring it. Therefore, the purpose is to move beyond the abstract boundaries of principled statements, striving instead to establish a framework that forms the logical foundation for properly managing the culture risk, which could aptly be described as the 'mother of all risks'. The stated insights may serve as a roadmap for risk managers who are tasked with addressing a significant and, in many respects, fundamental challenge. The remainder of this article, which is a theoretical paper based on conceptual analysis, is structured as follows: the first section explores definitions of risk culture and culture risk; the second outlines potential roles of corporate functions in mitigating culture risk. The third section examines the implications for the Risk Appetite Framework. The final section draws preliminary conclusions and sets the stage for future challenges.

## 1. Introduction

Financial institutions have long been tasked with managing risks intrinsic to their role in the financial markets. Risk management, therefore, represents a core component of financial intermediation (Allen and Santomero, 1997). As financial markets grow increasingly complex, the demands placed on risk management continue to expand, necessitating ever more sophisticated methodologies and techniques. As a result, risk management is characterized by high technical intensity, requiring the Chief Risk Officer (CRO) to possess a deep understanding of quantitative methods, data analytics, and other technical domains often rooted in the hard sciences. The CRO's role has become pivotal in not only monitoring existing risks but also anticipating and mitigating emerging threats. This strategic position demands a blend of technical proficiency, forward-thinking leadership, and the ability to collaborate across all levels of the organization to ensure the institution's resilience and compliance with evolving regulatory standards.

One of the primary implications of this evolution lies in the necessity of integrating the culture of control with a robust culture of risk. While it is widely recognized that risks and controls represent two sides of the same coin (Cocozza, 2024), emphasizing the 'culture of control' focuses on remedial interventions (addressing risks after they have materialized), while emphasizing the 'culture of risk' underscores the importance and prominence of preventative activities (mitigating risks before they arise). It goes without saying that both are essential; however, the latter may prove to be more effective and cost-efficient.

Therefore, technical expertise alone is no longer sufficient for risk managers in contemporary financial institutions. Their frequent involvement in strategic decision-making forums necessitates the development of additional competencies, including strong communication skills, the ability to influence stakeholders, and a nuanced understanding of organizational dynamics. These capabilities enable risk managers to effectively contribute to broader strategic discussions while ensuring that risk considerations are integrated into decision-making processes, with the ultimate scope of good governance, fundamental for the stability and safety of financial institutions, aligning with the overarching goals of Supervisors.

The Draft Guide on Governance and Risk Culture disseminated in July 2024 by the European Central Bank (ECB, 2024) emphasizes that shortcomings in risk culture can serve as early warnings for financial instability, making sound governance essential for strategic resilience and sustainable business operations. Looking ahead, the role of Risk Management appears to be taking on a central position within corporate dynamics. No longer viewed as a cost centre, it is emerging as a critical element in the value creation chain. Effective risk control necessitates a dual focus on both returns and risks by business line leaders, as well as the full engagement of senior management. In this perspective the culture risk, a topic fundamental to financial institutions (FIs) management, forces senior

leadership to set the tone for the organizational culture and possesses tangible and effective tools to do so, including incentive plans, capital allocation decisions, investments in control structures and resources, and the role assigned to the CRO and other control functions. Additionally, leadership must adopt a proactive and 'intrusive' approach to the decisions made by units responsible for assuming risk.

Within this framework, clearly defining what constitutes 'risk culture' distinguishing it from 'culture risk', determining the functional corporate actors of primary importance, outlining their potential roles and duties, and conceptualizing metrics for measuring culture risk remain topics that are, to some extent, yet to be fully explored. The concept of risk culture originates from an initial intervention by the Financial Stability Board (FSB, 2014), following the early regulatory initiatives on the Risk Appetite Framework (RAF). The main objective of the FSB (2014) was the development of a comprehensive framework for understanding and assessing risk culture within Fis. The document outlined the foundational elements of a sound risk culture, highlighting seminal themes that were subsequently revisited by various stakeholders. (BCBS, 2015; EBA, 2021; EBA 2025), particularly concerning the corporate governance of banks ([1]). Moreover, consistent risk culture accounts for Environmental, Social and Governance (ESG) risks implemented within the institution in accordance with EBA (2021). As this subject evolved, attention has progressively shifted to the governance at the upper echelons of organizations, the role of management body, and the significance of the RAF. In this context, both the role of corporate functions and the specification of a comprehensive set of culture risk indicators remain underexamined. The latter represents a main challenge.

The purpose of this paper, a theoretical paper based on conceptual analysis, is to address these questions, with the aim of bridging the gap between fundamental principles and the effective implementation of the intended objectives. The article is structured as follows: after defining the concepts of risk culture and culture risk through a deductive logical process (Section 2) and the addressing of relevant drivers (Section 3), the focus shifts to the roles that can be attributed to corporate functions and their potential responsibilities (Section 4), culminating in the identification of criteria useful for developing and maintaining culture risk indicators. (Section 5) The last section (Section 6) draws preliminary conclusions on the topic, while acknowledging that the subject is still under investigation, as it is significantly influenced, among other factors, by varying cultural perspectives that may emerge across different contexts.

## 2. Culture risk: a weighty challenge.

Culture, in its broadest sense, encompasses the collective values, beliefs, norms, and practices shared by a group of people. It serves as a framework that shapes behaviour, decision-making, and interactions within a social system. Culture is transmitted across generations through socialization and institutionalized practices, evolving over time as it adapts to environmental, historical, and societal changes. Anthropologically, culture is both material and symbolic, influencing tangible expressions (artifacts, systems) and intangible aspects (ideologies, shared meanings). It provides a cohesive identity to groups, guiding behaviour and ensuring continuity amidst diversity.

Accordingly, corporate culture refers to the specific set of shared values, norms, and practices that characterize an organization. It is both a product of and a contributor to the organization's identity, shaping how members interact internally and externally. Corporate culture influences decision-making, communication, and the prioritization of goals, aligning individual behaviours with organizational objectives. It arises from leadership philosophies, operational strategies, and the historical and social contexts within which the organization operates.

The key dimensions of corporate culture include values and norms, behavioural expectations, leadership and management styles as well as symbols and rituals (Figure 1)

| Values and Norms | Behavioural Expectations |
|---|---|
| • the foundational principles guiding actions and decisions | • unwritten rules and standards for conduct within the organization |
| **Corporate Culture** | |
| Leadership and Management Styles | Symbols and Rituals |
| • the tone set by leaders that shapes the organizational ethos | • practices and artifacts that reinforce a shared sense of purpose and identity |

*Figure 1: Dimensions of Corporate Culture*

Corporate culture plays a critical role in organizational performance, influencing innovation, employee engagement, adaptability to change, and ethical behaviour. Strong corporate cultures foster alignment between organizational goals and individual motivation, while fragmented cultures may lead to conflicts and inefficiencies. In summary, as in the popular quote apocryphally credited to management consultant Peter Drucker, 'culture eats strategy for breakfast', emphasising that a powerful and empowering culture is a sure route to success.

Risk culture can be regarded as a subset of the corporate culture. It specifically pertains to the norms, attitudes, and behaviours related to risk awareness, assessment, and management within an organization. It encompasses how risks are perceived,

---

[1] For further insights on the published works on the subject, the following are recommended: Bockius et al. (2024); Carretta et al. (2024); Kunz and Heitz (2021).

communicated, and addressed across all levels, influencing the organization's capacity to identify, mitigate, and respond to uncertainties.

Characteristics of risk culture include risk awareness, behavioural norms, communication practices, as well as accountability and incentives (Figure 2).

| Risk Awareness | Behavioural Norms |
|---|---|
| • the degree to which employees and decision-makers recognize and understand risks | • practices and attitudes toward risk-taking, caution, and accountability |
| **Risk Culture** | |
| Communication Practices | Accountability and Incentives |
| • mechanisms for reporting, discussing, and escalating risk-related concerns | • structures ensuring that individuals and teams are responsible for their roles in risk management |

*Figure 2: Characteristics of Risk Culture*

Given that banking activities, and financial intermediation more broadly, are inherently centred on risk – which, together with the financial resources collected, constitutes the core input of such activities – the presence of a robust risk culture is not merely desirable; it is a critical element actively pursued by supervisory authorities in their mission to ensure the ongoing safety and stability of banks. This necessity becomes even more pronounced in the current environment, where intermediaries face economic, competitive, and geopolitical challenges while simultaneously managing risks associated with climate change, environmental sustainability, and technological advancements. In fact, according to BCBS (2015, 2), recalling FSB (2014), risk culture is defined as «*a bank's norms, attitudes and behaviours related to risk awareness, risk-taking and risk management, and controls that shape decisions on risks. Risk culture influences the decisions of management and employees during the day-to-day activities and has an impact on the risks they assume*».

A strong risk culture aligns risk-taking behaviours with organizational objectives and regulatory expectations, fostering prudent decision-making and resilience. Conversely, a weak risk culture may result in misaligned incentives, insufficient risk controls, and an increased likelihood of operational or strategic failures. Leadership commitment, transparency, and continuous education are pivotal in embedding an effective risk culture within the broader corporate culture. Risk culture encompasses the collective mindset, norms, and behaviours that influence how risk is perceived, assessed, and managed within an organization. A strong risk culture aligns risk-taking conducts with organizational goals and regulatory expectations, fostering ethical decision-making and resilience.

Culture risk emerges when there is a divergence between the stated values of an organization and the actual practices and behaviours of its employees. This misalignment can lead to ethical lapses, operational inefficiencies, and reputational damage, ultimately compromising the institution's stability. Coherently culture risk may be referred to the potential adverse outcomes that arise from misalignments between an organization's stated values, norms, and principles and the actual behaviours, attitudes, and practices exhibited by its members. It encompasses risks stemming from deficiencies in fostering a cohesive and ethical culture that supports the organization's strategic objectives, regulatory compliance, and long-term sustainability.

Culture risk in financial institutions (FIs) can manifest in various forms, including ethical misconduct, operational inefficiencies, resistance to change as well as inadequate risk awareness, i.e. insufficient integration of risk management principles within the organizational culture, leading to poor decision-making or excessive risk-taking (Figure 3).

| Ethical Misconduct | Operational Inefficiencies |
|---|---|
| • failures in promoting integrity and adherence to ethical standards, potentially leading to reputational damage, legal violations, or regulatory penalties | • lack of alignment between cultural norms and operational practices, resulting in inconsistencies, inefficiencies, and errors |
| **Culture Risk** | |
| Resistance to Change | Inadequate Risk Awareness |
| • inflexibility or inertia in cultural attitudes that hinder adaptability to evolving market conditions, technologies, or regulatory requirements | • insufficient integration of risk management principles within the organizational culture, leading to poor decision-making or excessive risk-taking |

*Figure 3: Culture Risk Instances*

Supervisory authorities increasingly emphasize the management of culture risk as a fundamental component of FIs organizational governance, recognizing its critical role in mitigating broader operational, financial, and reputational risks.

Addressing culture risk requires ongoing leadership commitment, clear communication of values, and mechanisms for monitoring and reinforcing desired behaviours throughout the organization (Figure 4). In this respect, three fundamental pillars of culture risk management emerge as the foundation of the mechanism illustrated in Figure 4.

These pillars are:

- the leadership role, which must demonstrate profound risk awareness and a corresponding strong commitment, by communicating expectations clearly and consistently to reinforce a risk-aware culture;

- the effective communication throughout all levels of the organization, both top-down and bottom-up;

- the critical role of the organizational function, in addition to the corporate control functions.

Regarding the first pillar, the primary actors involved are the board of directors, board-level committees, and, where applicable, delegated executives. For the second pillar, it is essential to implement not only active speaking but also active listening (Cocozza, 2025). Finally, the third pillar requires the 'full maturity' of the organizational function, which serves as the primary safeguard of accountability, in conjunction with the human resources function (HR) for both incentives and induction and training programs.



*Figure 4: Foundational elements for addressing culture risk.*

Therefore, the palindrome 'risk culture & culture risk' is not an elegant pun.

Risk culture refers to the shared values, attitudes, and practices regarding risk awareness and management within an institution. Culture risk, on the other hand, arises from misalignments between an organization's stated values and the behaviours exhibited by its employees. The lack of a robust risk culture gives rise to culture risk, which can prove to be detrimental or even fatal to bank's stability and sustainability.

## 3. Culture risk – and value – drivers.

The aforementioned lack of a robust risk culture becomes a risk factor that impacts corporate performance in complex and multifaceted ways, many of which are not easily quantifiable in terms of their effect. Hence, promoting a robust risk culture serves as a comprehensive preventive measure against culture risk and, as such, it constitutes a fundamental component of the culture risk management process.

Central to this preventive framework is the concept of the 'tone from the top', which encompasses the ethical climate, cultural values, and behavioural standards set by an organization's senior leadership. The latter includes the board of directors, executive team, and other high-ranking officials. The tone from the top reflects the attitudes, decisions, and actions of senior leaders, demonstrating their unwavering commitment to organizational values, effective governance, and sound risk management practices.

Consistently, according to foundational elements reported in Figure 4, three warning signs can be immediately identified: the lack of independence, signalling insufficient commitment from leadership; the absence of adequate whistleblowing mechanisms, indicative of ineffective communication; and weak accountability, reflecting deficiencies within the organizational lines.

Indeed, these main red flags can be immediately identified for risk culture shortcomings. By addressing these issues, institutions can build resilient frameworks capable of adapting to evolving risks.

According to the ECB (2024, 12), as reported in Figure 5, risk culture components include, apart from the already mentioned 'tone from the top', effective communication challenge and diversity, incentives and accountability for risks. Root causes of culture risk and are identified as «*cultural drivers*».

- Composition of the management body and senior management
- Mandate of the management body to oversee the group
- Role of the management body in monitoring and assessing the risk culture

- Board's oversight and constructive challenge
- Risk-based strategic decisions
- Communication from the top on risk and compliance matters and core values to promote good behaviours among staff
- Reflection of the board and senior management on their own behaviour as well as actions as role models
- Board facilitates staff training in areas such as psychological safety and the bank's speak-up policy
- Reaction to supervision, in particular follow-up of findings related to risks and controls
- Role of leaders in fostering constructive and diverse decision-making; create group norm in which it is acceptable to challenge

- Connection of RAF to strategic processes
- Clear link between variable remuneration framework and risk appetite
- Clear documentation on the variable remuneration framework
- Long-term incentives considered in the remuneration and promotion framework
- Consequence management framework for misconduct (disciplinary processes and sanctions)

- Balance between risk and reward in daily decisions
- Appropriate incentives, including ex ante and ex post incentives on remuneration
- Balance of financial and non-financial performance criteria
- Applied metrics and limits commensurate with actual level of risk and risk appetite
- Transparency in promotion process and alignment with ethics. Misconduct reflected in promotion process

**Tone from the top and leadership**   **Incentives**

- Time commitment of management body members
- Escalation processes e.g. escalation and alert mechanisms for risk and control issues and findings
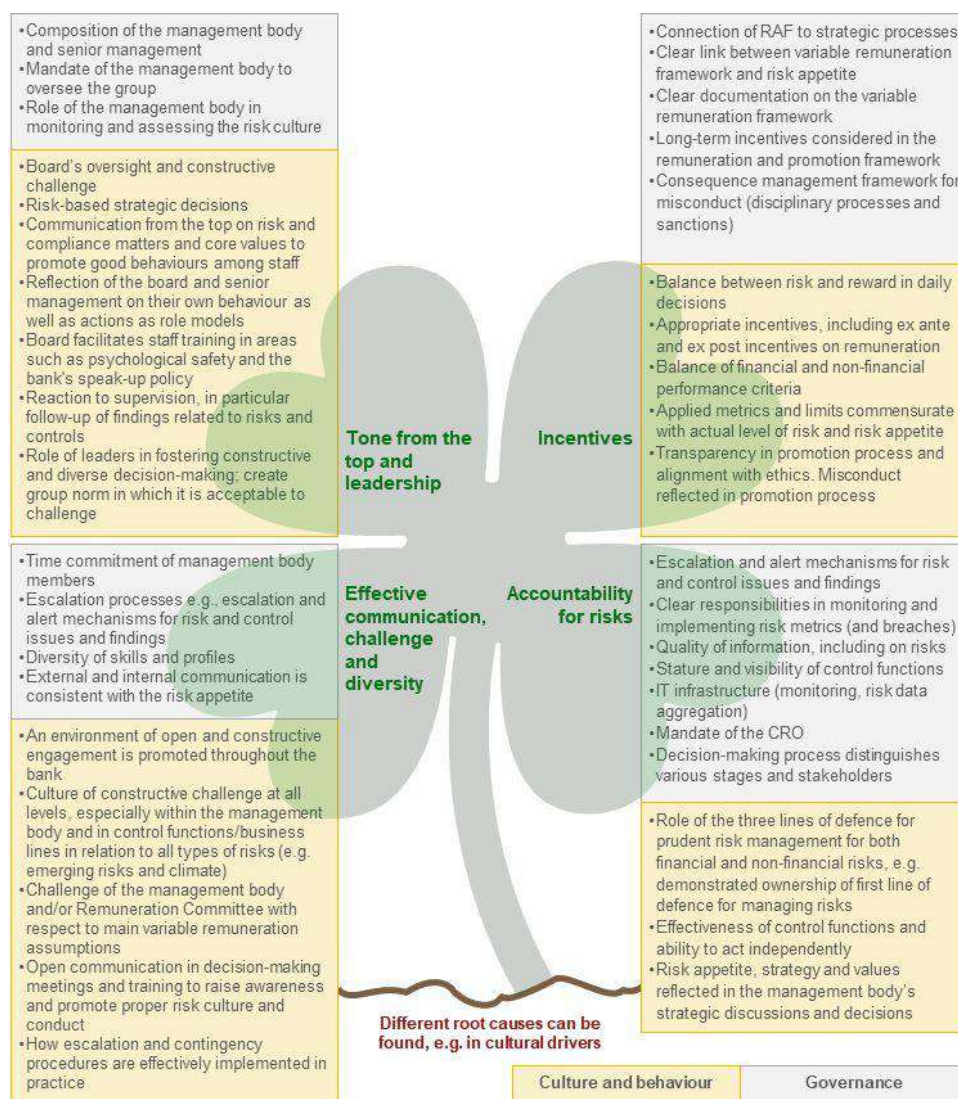- Diversity of skills and profiles
- External and internal communication is consistent with the risk appetite

- An environment of open and constructive engagement is promoted throughout the bank
- Culture of constructive challenge at all levels, especially within the management body and in control functions/business lines in relation to all types of risks (e.g. emerging risks and climate)
- Challenge of the management body and/or Remuneration Committee with respect to main variable remuneration assumptions
- Open communication in decision-making meetings and training to raise awareness and promote proper risk culture and conduct
- How escalation and contingency procedures are effectively implemented in practice

**Effective communication, challenge and diversity**   **Accountability for risks**

- Escalation and alert mechanisms for risk and control issues and findings
- Clear responsibilities in monitoring and implementing risk metrics (and breaches)
- Quality of information, including on risks
- Stature and visibility of control functions
- IT infrastructure (monitoring, risk data aggregation)
- Mandate of the CRO
- Decision-making process distinguishes various stages and stakeholders

- Role of the three lines of defence for prudent risk management for both financial and non-financial risks, e.g. demonstrated ownership of first line of defence for managing risks
- Effectiveness of control functions and ability to act independently
- Risk appetite, strategy and values reflected in the management body's strategic discussions and decisions

**Different root causes can be found, e.g. in cultural drivers**

| Culture and behaviour | Governance |

*Figure 5: Map of risk culture components, connecting governance, culture and behaviour. Source: ECB (2024, 12).*

As can be inferred from Figure 5, the ECB approach sets two main risk drivers for culture risk: 'governance' and 'culture and behaviour'. Consistently, the ECB (2024, 15) lists 'governance red flags' and 'behavioural and cultural red flags'. The occurrence of listed red flags, qualified «*non-exhaustive*», gives evidence of an inadequate risk culture and raises culture risk, as shown in Figure 6 and detailed in Figure 7.



Governance red flags + Behavioural and cultural red flags = Inadequate risk culture → Culture risk

*Figure 6: Culture risk drivers.*

The distinction between governance and behavioural and cultural risk drivers allows for further consideration on the individual cultural profiles of specific FIs. Smaller entities may be more susceptible to culture risk due to the traits of their governance structures and the impact of cultural and behavioural stratifications. Consequently, an inversion of the principle of proportionality may arise, suggesting heightened attention and caution, particularly in Less Significant Institutions (LSI). Therefore, it should in no way be construed as a plea for mitigating judgment. As a matter of fact, given the heightened complexity of implementing dedicated dashboards, LSIs must maintain a robust adoption process that is stringent yet free from undue bureaucratic encumbrances. Similarly, the adoption of a specific business model or a particular operational focus may serve as additional elements of individual profiling in relation to both control points and organizational structures. Similarly specific localization of activities, the prominence of cross-border operations, and the predominance of credit transactions characterized by peculiar contours can serve as significant factors in the customization of cultural risk profiles.

The proper conceptual framing of culture risk permits its classification within the category of governance risks, as defined by the EBA (2025). Culture risk can be conceptually situated within the broader category of governance risks – i.e., the 'G' component of ESG. The EBA explicitly define governance risks as encompassing deficiencies in executive leadership, ethical standards, and management practices, all of which can generate material financial risks that institutions are required to assess and manage. Although the term 'culture risk' is not explicitly employed, its conceptual features are clearly embedded within the EBA's treatment of governance: inadequate internal controls, insufficient board oversight, and failures in leadership behaviour are all identified as sources of governance-related vulnerabilities. The EBA further call for institutions to integrate ESG considerations into their standard risk management frameworks, emphasizing the role of ESG risks as potential amplifiers of traditional financial risk categories such as reputational, operational, and business model risks. In this context, the promotion of a sound risk culture – defined by effective communication, shared risk awareness, and clear accountability across all organizational levels – is deemed essential. In sum, the EBA's framework supports the interpretation of culture risk as an integral component of governance risk, insofar as it reflects the behavioural and ethical dynamics underpinning effective ESG risk management. This aligns with emerging technical literature (AIFIRM, 2025), which frames culture risk as a function of misalignment between formalized values and actual behaviours, highlighting its systemic implications for internal governance and institutional resilience.

Moreover, the capability to fully grasp risk culture shortcomings may also be shaped by the relevance the organisation and the board assign to 'risk and control' dimensions with respect to commercial aims. Although the academic perspective emphasises the equal importance of both, in corporate practice recognition of this equivalence sometimes encounters resistance. This may stem from differences arising from the depth and nature of experience in FIs management, as well as from the cultural and generational backgrounds of the individuals concerned.

| Risk culture dimension | Governance red flags | Behavioural and cultural red flags |
|---|---|---|
| Tone from the top and leadership | - Insufficient management body oversight of internal control functions and the management body in its management function<br>- Low number of formally independent members<br>- Insufficient subsidiary oversight<br>- Inadequate escalation and consequence management framework in the case of risk, ethical or compliance issues<br>- Inadequate conflict of interest policy and ethics framework | - Insufficient ownership of and responsibility for conduct risk<br>- Unsatisfactory tone from the top from the management body to promote good behaviours among staff<br>- Dismissive attitude among staff towards compliance, regulation and supervision<br>- Inadequate tone from the top on the balance of risks and rewards<br>- Concentration of power in a few members of the management<br>- Unethical behaviours not sufficiently sanctioned by the bank and insufficient communication on these issues |
| Culture of effective communication and challenge and diversity | - Deficiencies in the whistleblowing process<br>- Governance arrangements, including, committee structure and escalation process not facilitating debate<br>- Inadequate diversity framework | - Lack of challenge and debate within the management body (discussion dominated by a few management body members)<br>- Insufficient challenge of the management body in its supervisory function and/or its committees (e.g. remuneration committee) with respect to the main variable remuneration assumptions<br>- Insufficient challenge from internal control functions (e.g. lack of a role for the risk management function or its head in challenging decisions)<br>- Insufficient independence of internal control functions from the management body in its management function (e.g. filtering or review of information included in internal control function reports prior to the approval process)<br>- A culture of fear leading to an unwillingness to report mistakes, risk breaches or material concerns<br>- Lack of diversity (skills, gender, background) or inclusion, possibly contributing to "groupthink"<br>- Lack of meetings and training to raise awareness and promote proper risk culture and conduct |
| Incentives | - Documentation underpinning the variable remuneration framework (e.g. KPIs) either missing or ambiguously worded<br>- Lack of interplay between strategy and risk appetite<br>- Framework to address behaviours not aligned with prudent risk-taking<br>- Lack of link between variable remuneration framework and risk appetite<br>- Impaired consequence management (e.g. malus and clawback clauses exist only as a formality)<br>- Lack of individual accountability, including in the bank's remuneration and/or consequence management framework | - Incentive system does not incentivise desired behaviours<br>- Promotion process does not reflect conduct/misconduct, ethics and behaviour<br>- Applied metrics and limits are not commensurate with the bank's actual level of risk and its risk appetite<br>- Imbalanced deployment of financial performance criteria versus non-financial criteria<br>- Wrong incentives, e.g. remuneration of the CRO linked predominately to commercial objectives or connected with the performance of activities that the risk management function monitors |
| Accountability | - Low stature and understaffing of internal control functions<br>- RAF not comprehensive or well implemented<br>- Weak information technology (IT) and data aggregation framework<br>- Lack of a comprehensive "lessons learned" process to identify and address similar risks | - Unbalanced application of the third line of defence, i.e. the first line of defence lacking a culture of accountability for risk, leaving this to the second and third lines of defence<br>- Insufficient transparency in reporting (especially in the case of issues/concerns)<br>- Risk management seen as a barrier to achieving business objectives |

*Figure 7: Risk culture red flags (non-exhaustive list). Source: ECB (2024, 15).*

To provide practical application to the logical framework analysed here, it is appropriate to offer some concrete examples. The transition from theoretical postulation to practical applicability in the domain of culture risk necessitates a critical clarification concerning the measurability of the relevant constructs. In line with the managerial axiom that "what cannot be measured cannot be managed", it is imperative to delineate the object of measurement with precision. Specifically, it is not the degree or diffusion of risk culture per se that must be quantified, but rather the risk arising from its deficiency, i.e., culture risk. As previously argued, cultural inadequacies within an organization function as latent drivers of culture risk, which, like other risk categories, ultimately materialize through economic consequences, such as increased operational costs or reduced revenues. Accordingly, the operationalization of culture risk management must be anchored in the identification and deployment of appropriate Key Risk Indicators (KRIs), capable of capturing deviations from expected cultural norms and signalling potential misalignments before they escalate into broader governance failures.

Figure 8 presents, without claiming to exhaust the subject, a selection of processes that can be activated for the addressing of culture risk. Depending on the context, the outlined processes aim to establish an appropriate cultural climate, assess the current state within the individual organization, activate preventive mechanisms to mitigate exposure to culture risk, as well as implement traditional processes for identifying Key Risk Indicators (KRIs) and their corresponding monitoring.

| Leadership Commitment (Tone from the Top) | •Action: establish and demonstrate clear ethical standards, cultural values, and commitment to risk awareness. <br> •Impact: sets the foundation for a risk-aware culture through leadership example and strategic alignment. |
| --- | --- |
| Governance and Framework Development | •Action: develop policies, codes of conduct, and frameworks that integrate culture risk into governance structures. <br> •Impact: ensures culture risk is formalized and embedded into organizational processes. |
| Training and Awareness | •Action: design and implement training programs to educate employees on cultural expectations, risk management principles, and ethical behavior. <br> •Impact: builds capacity and awareness, empowering employees to contribute to a strong risk culture. |
| Behavioral Monitoring | •Action: use surveys, interviews, and performance reviews to assess employee behaviors, attitudes, and alignment with cultural values. <br> •Impact: provides insights into potential cultural misalignments and areas for improvement. |
| Preventive Mechanisms | •Action: implement systems for whistleblowing, escalation, and conflict-of-interest management to address potential issues early. <br> •Impact: mitigates culture risk proactively before it escalates into significant problems. |
| Continuous Feedback and Improvement | •Action: regularly review cultural practices, update frameworks, and act on lessons learned from audits and incidents. <br> •Impact: promotes adaptability and ensures the risk culture evolves with the organization |

*Figure 8: Processes for addressing culture risk.*

## 4. Culture risk: who is responsible for what?

Once the culture risk has been identified, it is appropriate, following the logical process typically adopted for other risk categories, to try to identify responsibility for its mitigation and promotion within corporate functions (Section 1).

As far as the mitigation is concerned, the responsibility falls certainly– although not exclusively – within the domain of the Internal Control Framework (ICF). The ICF is a cornerstone of governance in banking institutions, serving to ensure compliance, manage risks, and safeguard organizational integrity. Beyond its operational mandates, the internal control system is pivotal in fostering a robust risk culture and addressing culture risk. As FIs face increasingly complex challenges – including regulatory scrutiny, technological disruptions, and ESG concerns – the internal control system must evolve to address these demands, including also behavioural dimensions. The internal control system appraises culture risks through audits, behavioural assessments, and whistleblowing mechanisms. Clear escalation protocols, effective conflict-of-interest policies, and active reporting mechanisms address cultural misalignments proactively. The alignment of culture risk with governance requires that culture risk is integrated into the institution's governance frameworks, including the RAF, to ensure systematic management. Hence, regular evaluations of cultural practices and lessons learned from incidents strengthen the institution's ability to mitigate culture risk effectively. Addressing these aspects requires a comprehensive approach where the control function of second and third level play distinct yet interrelated roles. The compliance function ensures adherence to regulatory requirements, ethical standards, and internal policies. By assessing codes of conduct, providing training, and monitoring conducts, the compliance function shapes the ethical foundation of the organization. It also identifies and addresses misalignments that contribute to culture risk, fostering an environment where employees understand and embrace risk-aware practices. Risk management identifies, assesses, monitors, and mitigates risks that could impact the institution's objectives. Beyond managing traditional risk categories, the function integrates culture risk into the RAF and broader governance structures. By promoting proactive risk awareness and embedding accountability, risk management strengthens the organization's capacity to address cultural challenges. Internal audit provides independent assurance on the effectiveness of the organization's governance, risk management, and control systems. It assesses whether risk culture is embedded across the institution and identifies gaps in cultural alignment. Internal audit also evaluates the effectiveness of measures taken to mitigate culture risk, ensuring continuous improvement and accountability.

The timing of actions executed by internal control functions – compliance, risk management, and internal audit – is strategically aligned with the stages of risk and control activity within an organization. These stages, delineated as *ex-ante* (preventive), real-time, and *ex-post*, define the distinct responsibilities and levels of engagement for each function. In the preventive stage, the objective is to anticipate and mitigate risks before they materialize, thereby reducing the likelihood of adverse events. At this stage, the compliance function plays a pivotal role, enforcing regulatory requirements, organizational policies, and ethical standards designed to proactively address potential risks. Concurrently, risk management contributes by identifying emerging risks, assessing their potential impacts, and defining risk limits within the institution's RAF. Internal audit, however, typically has minimal involvement at this stage, as its primary responsibility is to deliver retrospective evaluations and assurance. The real-time stage focuses on the active monitoring and management of risks as they arise, ensuring timely and effective responses to mitigate potential impacts. During this phase, risk management assumes a leading role, continuously monitoring risk exposures, maintaining alignment with established thresholds, and making necessary real-time adjustments. The compliance function supports these efforts by ensuring ongoing adherence to regulatory and organizational standards amidst dynamic operations. Internal audit, though less central, offers moderate involvement by providing immediate feedback on control effectiveness and participating in oversight where required. In the *ex-post* stage, the emphasis shifts to the thorough analysis, evaluation, and enhancement of processes and controls following the occurrence of a risk event or control failure. At this juncture, internal audit assumes a dominant role, conducting in-depth investigations to uncover root causes, recommending corrective measures, and driving initiatives to strengthen organizational resilience. Risk management evaluates the broader implications of the incident on the institution's risk framework, revising mitigation strategies as necessary. Simultaneously, the compliance function ensures that any regulatory violations are properly identified, reported to relevant authorities, and addressed through corrective actions. This systematic alignment of functional roles across the risk management continuum ensures a coordinated, effective approach to risk governance, fostering organizational resilience, regulatory compliance, and strategic alignment. By optimizing the interplay of these functions, institutions can establish a proactive, agile, and comprehensive framework for managing risks across all phases of their operations.

The temporal distribution of responsibilities highlights the interdependence of internal control functions, and their complementary contributions ensure a cohesive approach to promoting risk culture and mitigating culture risk. These functions collaborate to align insights, strategies, and actions, creating a unified framework for cultural resilience. The timing of action for internal control functions underscores the strategic alignment of their roles in addressing risks across all stages of organizational activity. This framework establishes a robust foundation for advancing both academic inquiry and practical innovation in understanding the interplay between control functions and risk management settings and decisions within modern organizational contexts. It further underscores the critical importance of fostering a cohesive and integrated approach to strengthening risk culture while proactively addressing and mitigating cultural vulnerabilities. Therefore, the scope of activities attributable to the control functions is contingent upon the extent of risk culture dissemination within the organization and, consequently, the organization's level of awareness regarding culture risk related matters.

During the development phase, the promotion of risk culture is paramount. The internal control system strengthens the role of leadership in establishing an ethical tone, strengthening behaviours that prioritize risk awareness and regulatory compliance. In the foundational phase, comprehensive training programs designed by compliance and risk management functions serve to educate employees on risk management principles and cultural expectations, thereby advancing training and capacity building. In the maturity phase, internal control functions actively monitor employee behaviours and attitudes, offering feedback and recommendations to ensure alignment with institutional values, thereby facilitating behavioural monitoring and feedback. Once pervasiveness is achieved, the internal control system embeds risk culture into governance frameworks and operational processes, ensuring consistency and accountability across all organizational levels and fully integrating risk culture into governance.

With this respect, the operating area of a bank, encompassing various functions such as operations, technology, and back-office support, plays a focal role in embedding and sustaining a robust risk culture. As the operational backbone, this area ensures the institution's strategic objectives are translated into day-to-day activities while mitigating culture risk. Its responsibilities extend beyond traditional operations to fostering accountability, promoting transparency, and aligning operational practices with the institution's risk culture. The enhancement of accountability begins with the establishment of clear roles and responsibilities for all operational staff, ensuring that individuals are fully aware of their contributions to risk management and cultural alignment. Furthermore, structured escalation mechanisms are implemented to identify and address risk incidents effectively, ensuring prompt resolution and minimizing potential impacts. The integration of risk culture into operational processes involves translating organizational policies and strategic goals into actionable procedures, fostering consistency in decision-making and adherence to regulatory and ethical standards as well as ensuring that risk management becomes an integral part of all operational decisions. Strengthening communication channels is another essential responsibility of the operating area. By developing transparent and efficient frameworks for the exchange of information, the operating area facilitates the flow of critical risk-related insights across organizational levels. This open communication environment not only enhances collaboration between operational staff and control functions but also encourages constructive dialogue and the escalation of concerns. Additionally, the operating area supports the establishment and utilization of whistleblowing mechanisms, ensuring that employees can report unethical behaviour or cultural misalignments in a secure and confidential manner. In promoting ethical practices, the operating area ensures that the leadership's commitment to fostering a strong risk culture is translated into tangible actions throughout the organization. Operational processes are designed to reflect and reinforce the institution's core values, aligning operational goals with ethical standards and strategic objectives. Tools and training programs are developed to support employees in making ethical decisions, particularly in complex scenarios where risks must be carefully balanced against opportunities. Finally, the operating area is instrumental in aligning incentives with the organization's risk culture. By integrating adherence to risk culture and operational discipline into performance evaluations, the operating area ensures that employees are rewarded for behaviours that align with the institution's ethical and risk management standards. Compensation and promotion frameworks incorporate risk-awareness metrics, creating incentives that prioritize long-term organizational success over short-term gains. Employees demonstrating exemplary alignment with the institution's risk culture and ethical values are recognized and rewarded, further reinforcing the importance of cultural adherence.

Through these responsibilities, the operating area not only supports the operationalization of the institution's risk culture but also acts as a vital conduit for embedding ethical and risk-aware practices at every level of the organization. Its efforts contribute to a cohesive and resilient organizational environment where cultural and risk management objectives are seamlessly integrated into operational realities.

## 5. Culture risk: KRIs and KPIs.

A proper risk-culture framework can be established and survive only if it is made visible through a disciplined cycle of measurement and reporting; otherwise, it remains an abstract corporate mantra. For this reason, the starting point for any CRO is to weave culture-related Key Performance Indicators (KPIs) and Key Risk Indicators (KRIs) into the bank's planning architecture and business-model definition, strategic targets, Internal Capital Adequacy Assessment Process (ICAAP) and the Internal Liquidity Adequacy Assessment Process (ILAAP) risk-profile analysis and, finally, the RAF. The resulting 'KPI/KRI management system' must be strictly aligned with strategic objectives so that cultural ambitions are translated into operational thresholds: if a KPI expresses what the institution wishes to achieve, the corresponding KRI signals how far the underlying cultural drivers may endanger that objective. In a strong risk culture, this alignment enables proportional graduation of risk, prioritisation of monitoring effort and corrective action, and, crucially, creates an auditable bridge between tone-from-the-top statements and day-to-day behaviour. Supervisors now expect such traceability; the ECB (2024) explicitly stresses that culture must be «measurable, verifiable and proportionate to size, complexity and business model», and regards the adoption of an integrated KPI/KRI dashboard as evidence of its verifiability. Hence, establishing clear culture metrics is no longer a voluntary exercise but part of the prudential perimeter.

Regulatory bodies have consistently advocated the application of clear KPIs for the assessment of individual managerial performance. This paper expands the scope of the discussion by proposing a unified and coherent framework aimed at the identification and management of both KPIs and KRIs, as depicted in Figure 9.
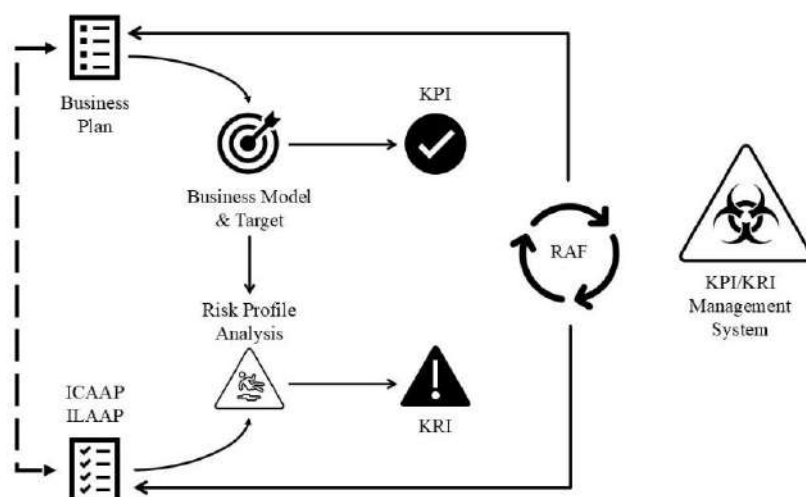


*Figure 9: KPI/KRI framework.*

Such a framework is posited as an essential component of robust corporate governance. Effective management necessitates rigorous planning, underpinned by an in-depth understanding of the organization, its operational mechanisms, target markets, and its exposure to potential risks, threats, and vulnerabilities. For the planning process to be truly comprehensive, it must also embody and communicate core values reflective of the organization's culture (Section 2).

The process of planning, defined as the articulation of business objectives within a specific business model, generates the essential data required to establish the institution's risk profile, as outlined in the ICAAP and the ILAAP. The effectiveness with which this process is executed constitutes the initial manifestation of an organization's risk culture. We emphasize the critical importance of a meaningful alignment between KRIs and KPIs, integrated within a coherent framework aligned with strategic goals. Such alignment is indispensable for the accurate classification and prioritization of risks, thereby enabling informed decisions regarding appropriate monitoring measures and potential mitigating actions. Obviously, KRIs aligned with KPIs must be pertinent (to the business model), measurable (in the most objective terms possible) and timely (reflecting environmental volatility and potential severity). In general, a sound KRI system is characterised by:

- details on the variables of the people-process-technology triad and on other corporate attributes most relevant to the proper functioning of the organisation in pursuit of strategic targets;

- classification of corporate assets according to their criticality for the bank;

- identification of the risks, threats and vulnerabilities the bank must face, based on probability of occurrence, operational and financial impact, and the organisation's capacity to mitigate the event;

- subsequent classification of risks, threats and vulnerabilities in terms of potential damage;

- linkage between key corporate objectives/KPIs and the most significant risks, in order to identify areas requiring enhanced monitoring and control;

- definition of parameters that determine when and how an identified risk becomes a serious threat;

- codification of a continuous process for reviewing KRIs and their metrics so as to detect any changes that require review and/or corrective action.

As any other risk management process, the identification phase is anchored to the people-process-technology triad and proceeds in four logical steps. First, the CRO maps critical assets – human capital, core processes, IT platforms – against their relevance to cultural objectives. Second, potential threats and vulnerabilities are classified by probability, severity and organisational ability to mitigate; this turns qualitative cultural ambitions into risk-sensitive categories. Third, each material risk is linked to the strategic KPI it could derail, thereby highlighting areas requiring heightened monitoring; and fourth, quantitative or qualitative parameters are set that allow objective detection of early deviation.

The resulting menu of KRIs will typically include

- financial indicators (e.g. risk-adjusted revenue versus conduct events);

- human-resources indicators (e.g. regretted turnover in control functions or training-completion ratios);

- operational indicators (e.g. process-break rates, near-misses, override frequencies);

- technology indicators (e.g. percentage of critical systems with end-of-life components);

- cyber-security indicators (e.g. phishing-simulation failure rates).

For each KRI the CRO establishes tolerance thresholds and escalation triggers. Because culture risks evolve with 'novel' external pressures – ESG litigation, AI bias, geopolitical disinformation – KRIs must be reviewed at least annually, with ad-hoc revisions whenever the business plan or the regulatory environment changes. Embedding this architecture in governance is the responsibility of the board and its committees. Under the assumption of 'full cultural maturity', the board approves the list of culture KPIs/KRIs as part of the RAF and receives regular dashboard reports, thereby making cultural performance a standing agenda item. In this process, Risk Management owns the design and validation of indicators, Compliance tests alignment with regulatory conduct expectations and Internal Audit provides ex-post assurance on data reliability and on the effectiveness of escalation.

Active leadership participation is indispensable: senior executives must use the dashboard in performance dialogues, and business-line heads must feel personal accountability for deviations; only then do KRIs become 'lived' metrics rather than compliance artefacts. It is also recommended to embed KPI/KRI attainment into variable-remuneration scorecards and claw-back clauses: for example, failure to close substantiated whistle-blower cases within target time would reduce the bonus of the manager concerned. Equally, exemplary adherence – such as proactive challenge of group-think – triggers positive recognition, converting cultural principles into tangible incentives.

Monitoring and escalation complete the cycle. A robust data-governance backbone feeds automated dashboards – accessible online and offline on mobile devices – to all three lines of defence, ensuring timeliness and transparency. Threshold breaches are colour-coded (e.g. amber for tolerance limits, red for breaches) and routed through predefined channels: first line rectifies operational issues; second line validates remediation and, if systemic, proposes RAF reviews; third line assesses lessons learned. Parallel whistle-blowing statistics, staff-survey sentiment scores and diversity metrics complement hard data, addressing ECB concerns that a 'culture of fear' could suppress early signals. Where red flags accumulate, for instance, weak management-body challenge, poor conflict-of-interest disclosure, or inadequate variable-pay documentation, the CRO must initiate a culture-risk incident report, triggering board scrutiny and, when necessary, supervisory notification. This closed-loop process proves vital in LSI, where proportionality must not reduce

vigilance; indeed, the inversion principle, as noted, warns that smaller entities may require greater indicator granularity, as informal governance can mask early cultural erosion.

The system described is complex and, to be robust and thus effective, must be based on a very solid process capable of managing the 'KRI life-cycle' (identification, assessment, monitoring, reporting to recipients). This calls for a precise allocation of responsibilities which, returning to Figure 9, could fall to the drafter of the RAF, allowing a centralised management of the implementation issues that arise in this field. Numerous points warrant attention. First, active participation by senior figures in the use of KRIs as an integral part of an enterprise risk-management programme must be ensured, without neglecting all other stakeholders in business and staff areas. Involving people – a qualitative building-block of an organisation's culture – facilitates the sharing of ideas and the use of indicators that are well understood by all. This requires identifying parameters that are «*measurable*» and «*comprehensible*»: dashboards are an effective means of presenting information and facilitate its use. A continuous activity must also be defined to monitor, measure and analyse any changes in the metrics. Finally, the system must ensure that actions are generated whenever deviations from KRI metrics occur. This represents the ultimate confirmation that the system is effectively in use and that risk culture is not an abstract element in corporate culture.

Regular monitoring of KRIs is essential, but the frequency depends on the nature of each indicator. Some KRIs may require daily attention, while others can be reviewed monthly or quarterly. It is vital to establish a routine that corresponds to the potential impact and probability of the risk, once again in relation to organisational complexity and business model.

In conclusion, designing an effective system of KRIs aligned with KPIs entails several challenges. The main one, in our view, is not to overlook the need to align KRIs with KPIs; conversely, the risk of a lack of responsiveness to periodic measurements must also be managed, often justified by the absence of thresholds. Naturally, the greater the organisational complexity, the more technological challenges arise: the possibility of working through dashboards is based on system integration, often hampered by obsolescence and fragmentation of systems. We can summarise the foregoing by recalling that 'corporate culture' cannot be separated from the 'culture of reporting', which covers the entire process from KPI identification and data processing to presentation and active use.

Admittedly, the embryonic stage of development in this field does not yet allow for a fully articulated framework that extends beyond the best practices outlined in the Guide (ECB, 2024), an observation extensively discussed in the dedicated AIFIRM Position Paper (2025), that we recall here for a more detailed examination. Building on the conceptual architecture developed in the present contribution, several lines of inquiry warrant further attention in order to advance the theoretical robustness and operational applicability of culture risk frameworks. First, the proposed alignment between KPIs and KRIs calls for empirical validation across diverse institutional contexts, with a view to assessing its predictive efficacy and practical enforceability. In-depth case studies of cultural failures in financial institutions, aimed at tracing observable misconduct or governance lapses back to early cultural warning signs and deficiencies in oversight mechanisms, could serve as a valuable support in the practical definition of indicators. At the organizational level, further research should interrogate the behavioural drivers of cultural integrity, including leadership tone, groupthink dynamics, psychological safety, and the structuring of incentives.

From a governance perspective, the integration of culture risk into the RAF requires greater technical articulation, particularly concerning the calibration of tolerance thresholds, the linkage to capital planning, and the definition of escalation protocols. In parallel, the ongoing digital transformation of financial services raises urgent questions about how culture risk manifests in algorithmic environments and technology-led business models, especially among fintech entities operating outside traditional governance structures. Finally, the incorporation of qualitative data streams, such as employee sentiment analysis, whistleblowing metrics, and behavioural surveys, into formal risk governance frameworks may offer a promising path toward the early detection of cultural vulnerabilities, provided that appropriate safeguards for data integrity, confidentiality, and accountability are in place. Collectively, these research directions offer a coherent agenda for advancing the state of the art in culture risk management, moving the field beyond principled aspirations toward verifiable, actionable, and institutionally embedded practices.

## 6. Conclusions

If culture risk constitutes – as actually is – a genuine risk category, it must be addressed on par with other risk types. Identification, measurement, and management should serve as the foundational elements for developing an effective mitigation of this specific risk category. In this way, the FI's strategy can be nourished and strengthened by a well-established risk culture, setting a way of creating organizations that are flexible and innovative and where individuals take responsibility for results – moving away from bureaucratic silos where formulaic approaches dominate. In other words, risk culture represents shared norms, attitudes, and behaviors toward risk management and awareness at all levels.

In the end, the 'risk manager of last resort' is the Chief Executive Officer (CEO), who bears ultimate responsibility for both results and risks. It is only from the top that a successful strategy and a corporate culture genuinely grounded in the understanding and control of risks can be effectively established and enforced. Nevrtheless, an effective culture-risk architecture hinges on a disciplined KPI/KRI ecosystem that transforms ethical aspirations into measurable managerial practice.

The translated framework underscores four imperatives:

1. embed KPI/KRI design in strategic planning so that risk-profile analysis (ICAAP/ILAAP) informs – and is informed by – the Risk Appetite Framework;

2. align every KRI with a corresponding KPI, thereby enabling graded risk prioritisation and proportionate corrective action;

3. maintain a dynamic life-cycle for indicators – definition, validation, monitoring, escalation – supported by clear ownership, dashboard-based transparency and thresholds that trigger timely intervention;

4. balance quantitative and qualitative signals across finance, human resources, operations, technology and cyber-security, with special vigilance for emerging 'novel risks'. When senior leadership visibly employs this dashboard in performance dialogues, the bank converts abstract cultural principles into operational discipline, ensuring that 'what gets measured gets managed' remains true even for the elusive domain of corporate culture.

Finally, implementation must respect proportionality while preserving comparability. Two design principles remain universal: bidirectional integration implying that KRIs flow into strategic KPI assessment, and KPI shifts trigger KRI re-validation; and actionability, that is to say every indicator must have an owner, a documented escalation path and a predetermined management response. When these principles are honoured, cultural-risk metrics cease to be a regulatory burden and become a strategic asset: they enable management to balance innovation and prudence, reassure supervisors, and, ultimately, protect stakeholder confidence in a volatile environment. As practitioners know, *what gets measured gets managed: culture risk is no exception*.

## References

AIA (2015). *La cultura del rischio*. https://www.aiiaweb.it/la-cultura-del-rischio.

AIFIRM (2024). *Comments to ECB Guide on governance and risk culture.* Available at: https://www.aifirm.it/wp-content/uploads/2024/10/2024-45-Risposta-cons.-BCE-Draft-Guide-on-Governance-and-Risk-Culture.pdf.

AIFIRM (2025). *Governance and Risk Culture*. Available at: https://www.aifirm.it/wp-content/uploads/2025/06/2025-Position-Paper-48-Governance-e-Risk-Culture.pdf.

Allen F. and Santomero A. (1997), The theory of financial intermediation, *Journal of Banking & Finance*, 21, pp. 1461-1485.

BCBS (2015). *Corporate governance principles for banks*. Available at: https://www.bis.org/bcbs/publ/d328.pdf.

Bockius H. and Nadine Gatzert N. (2024). Organizational risk culture: A literature review on dimensions, assessment, value relevance, and improvement levers, *European Management Journal*, 42 (4), pp. 539-564. https://doi.org/10.1016/j.emj.2023.02.002.

Carretta A., Fattobene L., Graziano E. A. and Schwizer P. (2024). Errors and Misbehaviors in banking and finance: a Systematic Literature Review and an Integrative Framework, *Journal of Management Governance*. https://doi.org/10.1007/s10997-024-09727-7. Available at: https://rdcu.be/d5VNQ.

Cocozza R. (2024). Fattori critici di successo del Risk Management: qualche istruzione per l'uso, *Rivista Bancaria Minerva Bancaria*, 3-4, pp. 57-84. https://dx.doi.org/10.57622/RB2024-03-C-04.

Cocozza R. (2024). Risk management, the board and the C-suite: The adaptive art of communication in times of change, *Journal of Risk Management in Financial Institutions*, 18, pp. 14-25. https://dx.doi.org/10.69554/nyap5618.

EBA (2021). *Guidelines on internal governance under Directive 2013/36/EU*. Available at: https://www.eba.europa.eu/sites/default/files/document_library/Publications/Guidelines/2021/1016721/Final%20report%20on%20Guidelines%20on%20internal%20governance%20under%20CRD.pdf.

EBA (2025). *Guidelines on the management of environmental, social and governance (ESG) risks*. Available at: https://www.eba.europa.eu/sites/default/files/2025-01/fb22982a-d69d-42cc-9d62-1023497ad58a/Final%20Guidelines%20on%20the%20management%20of%20ESG%20risks.pdf.

ECB (2016). *SSM supervisory statement on governance and risk appetite*. Available at: https://www.bankingsupervision.europa.eu/ecb/pub/pdf/ssm_supervisory_statement_on_governance_and_risk_appetite_201606.en.pdf.

ECB (2024). *Draft guide on governance and risk culture*. Available at: https://www.bankingsupervision.europa.eu/framework/legal-framework/public-consultations/pdf/ssm.pubcon202407_draftguide.en.pdf.

FSB (2014). *Guidance on Supervisory Interaction with Financial Institutions on Risk Culture*. Available at: https://www.fsb.org/uploads/140407.pdf.

Kunz J. and Heitz M. (2021). Banks' risk culture and management control systems: A systematic literature review. *Journal of Management Control*, 32, pp. 439-493. https://doi.org/10.1007/s00187-021-00325-4.

Romito F. (2012). L'evoluzione del Risk Management nelle banche: non solo misurazione, in P. Prandi (a cura di), Il Risk Management negli Istituti di Credito. Come affrontare le sfide in scenari di incertezza, FrancoAngeli, Milano pp.19-26.

# Integrating Machine Learning and Rule-Based Systems for Fraud Detection: A case study based on the Logic Learning Machine

**Pier Giuseppe Giribone (University of Genoa, BPER Group); Giorgio Mantero (Rulex Inc.); Marco Muselli (Rulex Inc.), Damiano Verda (Rulex Inc.)**

**Corresponding Author: Giorgio Mantero (giorgio.mantero@rulex.ai)**

## Abstract

Money laundering is one of the most relevant global challenges, with significant repercussions on the economy and international security. Identifying suspicious transactions is a key element in the fight against the phenomenon, but the task is extremely complex due to the constant evolution of the strategies adopted by criminals and the great amount of data to be analyzed daily. This study proposes a hybrid method that integrates Machine Learning models with heuristic rules, with the aim of identifying fraudulent transactions more effectively. The dataset used, SAML, includes millions of bank transactions and presents a strong imbalance between classes (fraudulent vs regular transactions). The entire process was carried out through a self-code platform designed to optimize data management, processing and analysis. The heuristic rules were evaluated using the covering and error metrics and then integrated into the Logic Learning Machine (LLM) task. The effectiveness of the approach was verified by comparing two main configurations: one based exclusively on the use of LLM and the other combining LLM and heuristic rules. The results obtained highlight that the integration of heuristic rules improves the performance of the model, confirming the synergy between Machine Learning and expert knowledge. This study confirms the effectiveness of the hybrid approach and emphasizes the importance of the union between automated analysis and human insight to address the challenges posed by money laundering.

**Key Words**: Anti-Money Laundering (AML), Transaction Monitoring, Synthetic Dataset, Machine Learning (ML), Heuristic Rules, Logic Learning Machine (LLM)

**JEL code:** C38, C45, K14, K22

## 1) Introduction

As described by the United Nations Office on Drugs and Crime (UNODC): "Money laundering is the processing of criminal proceeds to disguise their illegal origin. This process is of critical importance, as it enables the criminal to enjoy these profits without jeopardizing their source" (UNODC, 2021). Money laundering therefore indicates all those processes implemented by criminal organizations in order to disguise the origins of money obtained through illegal activities, such as corruption, drug trafficking, or fraud, to make it appear legitimate. By implementing it, such organizations can integrate illicit funds into the financial system, allowing further investment in illegal operations while still managing to avoid recognition by supervisory bodies.

These illicit activities pose serious concerns for the global economy because they are used to fund criminal activities, they disrupt financial markets and ultimately may also damage the reputation of the financial institutions involved in fraudulent activities.

Although it is clearly not possible to directly measure the extent of money laundering as we usually do with legitimate economic activities, its scale is massive, thus representing a significant threat to global financial systems. The UNODC estimates that money laundering accounts for 2-5% of global GDP annually, i.e. between 800 billion and 2 trillion EUR (UNODC, 2021), thus underscoring the need for robust detection and prevention mechanisms.

Money laundering is a global phenomenon, but data shows that it is more prevalent in specific industries and countries. In particular, the main sectors most affected by it are the real estate, the financial system, the gambling system, the trade of luxury goods, international trades, the construction sector and FinTech. Regarding the most problematic countries, the Financial Action Task Force (on Money Laundering), better known as FATF, identifies the countries with severe weaknesses in measures to combat money laundering and terrorist financing through the "High-risk jurisdictions subject to a call for action" list (blacklist) and the "Jurisdictions under increased monitoring" list (greylist).

Anti-Money Laundering (AML) thus refers to the regulations, policies, and procedures designed to detect and prevent money laundering activities. The main objective of AML agencies is indeed to prevent criminals from using the financial system to conceal the funds arising from their illicit activities. The work of AML professionals is divided into three main areas:

- **Prevention:** This part involves ensuring that governments, companies and financial institutions take all the necessary preventive measures to identify suspicious activities in their work;

- **Monitoring:** This consists of creating monitoring systems to inspect transactions and flag those that may be linked to laundering activities;

- **Reporting and prosecution:** This phase involves reporting any suspicious activity to Financial Intelligence Units (FIUs) in order to take action. In such sense, cooperation between bodies to enable criminal investigation is crucial.

All these procedures are carried out in full compliance with national and international laws (Financial Action Task Force, 2003). Despite advancements, AML efforts face several challenges:

- **Scalability:** The sheer volume of daily financial transactions demands highly efficient AML systems. For instance, international banks process millions of real-time transactions each day, requiring models that balance accuracy and speed for timely fraud detection without compromising system performance;

- **False positives:** AML detection systems often use strict rules and conservative thresholds to avoid missing fraud, but this leads to high false positive rates, driving up costs and causing delays that can harm customer relationships. The key is finding a balance between swiftness and accuracy, with systems that are fast yet precise enough to catch suspicious cases. The choice between highly accurate but slower systems and faster, less precise ones largely depends on the volume of alerts they need to handle;

- **Adaptability:** As criminal techniques evolve, and new fields emerge, such as cryptocurrencies and decentralized finance, AML systems must continuously adapt to keep pace. Relying only on past fraud patterns risks bringing rapid obsolescence, making it harder to detect new illicit behaviors. This requires constant updates to both heuristic rules and data-driven models, which can be costly, both in terms of maintaining high performance models and training operators with deep expertise.

The complex and multifaceted fight against money laundering is a problem that requires robust regulatory frameworks to ensure the resilience of the financial system. To counter the threats posed by it, governments, international organizations and financial institutions developed a wide range of guidelines and regulations over time with the aim of preventing, detecting, and prosecuting illicit activities. In this sense, part of our heuristic rules was drafted in order to align with these international guidelines and incorporate a regulatory perspective.

This study aims to investigate whether the combination of Machine Learning systems with heuristic rules can improve the effectiveness in detecting fraudulent transactions in the AML field. The proposed approach aims to exploit the strengths of data-driven methodologies while integrating specific sectorial expertise. The analysis was conducted using the Rulex Platform, an advanced platform that combines data analytics tools, Machine Learning tasks and heuristic rules management, allowing to efficiently implement and test models. The purpose of this study is to assess the effectiveness of this solution based on the Logic Learning Machine (LLM) algorithm by comparing it with other traditional Machine Learning approaches in order to analyze its performance in detecting odd activities in financial transactions datasets. We have chosen the LLM algorithm because it has proven to be valuable in solving problems in the context of financial and credit risk management. In particular, it has been employed both in an asset allocation context in order to select the optimal weights of a portfolio of ESG assets (Gaggero et al., 2024) and to improve models for predicting the probability of default in a set of U.S. companies (Berretta et al., 2025). In both contexts, it proved to be a reliable supervised Machine Learning technique characterized by a high level of explainability.

## 2) Literary review

The techniques used in the AML area mainly focus on the implementation of analytical and technological methodologies to detect and prevent money laundering activities. Specifically, there are customer due diligence measures (KYC and customer risk scoring), computing techniques to discover fraudulent patterns (ML algorithms) and finally manual investigation to confirm the flags raised. Traditionally, fraud detection has relied on deterministic rules, often derived from regulations, and subsequent manual checks by operators. These methods, albeit very useful, show severe limitations in terms of scalability and ability to adapt to complex and evolving fraudulent schemes, as well as high costs in terms of training suitable personnel.

In recent years, Machine Learning has emerged and begun to revolutionize fraud identification, allowing specialists to analyze large volumes of data and identify complex patterns that are not easily detectable with static rules (Teradata, 2022) (Nweze et al., 2024).

Machine Learning techniques are broadly categorized into supervised, unsupervised, and semi-supervised approaches, and are applied across several analytical dimensions, including anomaly detection, risk scoring, behavioral modelling and link analysis. Specifically, supervised models rely on labeled datasets distinguishing between normal and suspicious transactions. Notable algorithms include: Support Vector Machines (SVM), which are effective in high-dimensional spaces though computationally intensive on large, imbalanced datasets; decision trees, which are highly interpretable and useful for risk scoring and profiling but prone to overfitting if not appropriately pruned; and Radial Basis Function Networks (RBFN) which offer great adaptability and fast learning, but suffer from the risk of overfitting in cases of low feature diversity. On the contrary, unsupervised techniques cluster data without any prior label, making them especially suitable when suspicious labels are scarce. The most prominent algorithms are clustering techniques like K-means and CLOPE, which are used to group similar patterns of transactions to identify anomalies, and Expectation-Maximization (EM) methodologies to model customer behavior and detect deviations. Semi-supervised approaches are designed in such a way to strike a balance between the need for labeled data and the complexity of fraudulent patterns. The latter combine supervised learning and clustering, often using synthetic data to overcome the problem of class imbalance. Deep learning refers to a subclass of Machine Learning techniques that utilizes models composed of many layers of nonlinear transformations. These networks are capable of learning complex and abstract representations of data, often with performance far superior to other techniques, but at the cost of needing extensive computational resources and lacking interpretability. Another highly developed branch is the one related to graph-based methodologies and Social Network Analysis (SNA). The latter are increasingly used to model relationships among different entities, showing structural patterns for money laundering schemes. It works by building a graph where the nodes represent different entities (e.g. bank accounts) and the arcs represent the relationship between nodes (e.g. money transactions). The objective is to identify specific money laundering schemes like circular-shaped transactions or hub-and-spoke structures, characteristic of layering (Chen et al., 2018).

More recently, literature highlights a growing interest in hybrid methodologies that combine the rule-based approach with Machine Learning models. The latter are gaining more and more ground, as they leverage the strengths of both techniques. These approaches make it possible to improve the explainability of decisions, exploiting the expert knowledge embedded in heuristic rules, while maintaining the flexibility and learning ability of Machine Learning models.

Although literature has demonstrated that a rational integration of heuristic rules with Machine Learning models generally improves the fraud detection process, the creation, but above all the management, of large rulesets can lead to severe operational complexities

and difficulties in their interpretability. For this reason, in recent years research focused not only on the combination of the two methodologies, but also on the optimization of such sets of rules, with the aim of reducing their number and their computational cost, obviously without compromising the overall model performance. This particular field of research led to the development of models capable of improving sets of rules using techniques that, once again, combine algorithms and human expertise.

One of the most interesting studies in this field is that related to the development of the RUDOLF system (Milo et al., 2018). In this study the authors try to overcome the classic problem linked to the use of "mining" and Machine Learning techniques for the derivation of rules in the anti-fraud field. Since heuristic rules must necessarily be updated or redefined from time to time to keep up with fraud trends, researchers developed RUDOLF, a system that assists experts in defining and redefining the rules for identifying fraudulent transactions. RUDOLF first tries to uncover illicit instances by generalizing the initial rules proposed, and then it specializes them in order to avoid capturing unhelpful legitimate transactions. The changes are not mandatory but just proposed by the system: the supervisor can consequently accept, modify or reject each suggestion, based on his judgement and experience. This process goes on until the expert obtains the desired ruleset. Modifications to rules made by RUDOLF are associated with a cost-benefit model. This model assumes that every operation performed leads to a cost, but also to an entailed benefit, measured in terms of an increase in the number of frauds captured by the new rule or a decrease in the number of normal transactions captured. Therefore, the system's objective is to modify the existing ruleset so that the cost function is minimized.

In the performance comparison between the baseline version of RUDOLF and RUDOLF⁻, a variant that automatically refines the ruleset without consulting experts, the authors demonstrated that RUDOLF performs best in various domains, especially in the quality of predictions. This demonstrates the relevance of incorporating the expertise of AML domain specialists in the drafting of anti-fraud rules.

Another relevant contribution in this field is given by the ARMS project (Aparício et al., 2020). The Automated Rules Management System (ARMS) is a technique that optimizes the set of rules used thanks to heuristic search and a loss function defined by the user. Its proposed goal is to minimize the number of rules and alerts, while preserving the initial performance.

Instead of considering and evaluating each rule independently, researchers built a system to manage rules that considers the interactions between rules with different actions and priorities. The latter are needed because transactions may trigger different rules with contradictory actions, creating the need for a stable hierarchy between rules. The ARMS optimization process is basically implemented in two ways, specifically by disabling inefficient rules and by changing rules' priorities. The heuristic methods tested by authors are random search, greedy expansion and genetic programming. The results obtained on two big online datasets show that ARMS was able to remove almost 50% (and 80% in the second case) of initial rules, while maintaining the original performance of the system.

## 3) Regulatory Framework

The complex and multifaceted fight against money laundering is a problem that requires robust regulatory frameworks to ensure the resilience of the financial system. To counter the threats posed by it, governments, international organizations and financial institutions developed over time a wide range of guidelines and regulations with the aim preventing, detecting, and prosecuting illicit activities.

We now briefly explore the key international standards and regulatory frameworks that underpin AML efforts. These regulations provide a comprehensive framework for addressing financial crimes by establishing obligations for financial institutions and governments, promoting cross-border cooperation and ensuring rule-compliance through rigorous monitoring mechanisms.

## 3.1) FATF Recommendations (Financial Action Task Force)

Less than a year after its establishment by the G7 summit, the FATF issued in 1990 a report containing the well-known set of "Forty recommendations", later revised in 1996. Five years later, in 2001, the FATF expanded its mandate to also cope with the problem of terrorism financing (AML/CFT Anti-money laundering and Combating the Financing of Terrorism) and it continued to update its agenda over the years.

With its recommendations, the FATF is internationally endorsed as the global standard against money laundering and terrorist financing, as well as the financing of proliferation of weapons of mass destruction.

The recommendations, last comprehensively revised in 2012 and subsequently updated on specific issues such as virtual assets and beneficial ownership (e.g., 2021), set the minimum standards that countries must implement according to their specific circumstances and regulatory system (FATF, 2025). The FATF provides guidance on the following topics:

**Risk-based approach:** This means that each country should assess the risks that it faces and take appropriate preventive action in response (KYC-Chain, 2020). Such an approach may also be "scalable", in the sense that riskier instances clearly need more stringent measures and vice versa, so it is a proportionate approach.

**Sanctions:** The FATF recommends applying member states to implement a "targeted financial sanctions regime" to fully comply with the United Nations Security Council Resolutions (UNSCRs). The latter asks governments to freeze without delay the assets and funds of listed people or entities (or groups of them) that pose terrorist financing risk, and also ensure that no further financial assets are made available to them in the future (KYC-Chain, 2020) (FATF, 2013).

**Customer Due Diligence and record-keeping:** This principle states that financial institutions should not keep anonymous accounts or with obviously fictitious names. Additionally, it requires financial institutions to undertake customer due diligence measures, like identifying and verifying the identity of their clients (FATF, 2025). Furthermore, agencies are asked to keep records of all relevant information about clients in order to assess the risks posed by potential and current customers.

**Reporting of suspicious transactions and compliance Reporting:** it is a vital tool for AML and CFT measures to work in an efficient way. If authorities do not receive any reports, it is obvious that finding illicit activities becomes problematic. The FATF strongly

recommends that institutions implement a mandatory reporting obligation, regardless of the gravity of the illicit action, also posing great importance on the celerity in sending such reports: the sooner the better (KYC-Chain, 2020) (FATF, 2025).

**New technologies:** The FATF strongly suggests countries to be aware of how fraudsters may use new arising and disruptive technologies in order to commit crimes, particularly regarding the financial sector. In this sense institutions should not release new products or technological developments unless a prior risk assessment has occurred. This means implementing a sophisticated system to prevent, or at least better manage, any potential risk that could emerge. Regarding virtual assets (e.g. cryptocurrencies), countries should ensure that the service providers are regulated for AML/CFT purposes, registered and adequately monitored (KYC-Chain, 2020).

## 3.2) European Union Anti-Money Laundering Directives (AMLD)

The evolution of anti-money laundering efforts at European level can be described by analyzing the series of Directives produced. The first commitment dates back to 1991, when the first Directive was adopted to prevent the misuse of the financial system for the purpose of money laundering (European Commission, 2023). This process was undoubtedly influenced by FATF recommendations produced on the same year and driven by the rising international awareness and effort in tackling money laundering. In particular the first AMLD focused on the imposition of due diligence measures for financial institutions and on the establishment of a reporting system for all suspicious transactions. In the following Directives, the European Union kept revising the previous regulatory limitations and expanding the operational framework in order to mitigate the risks related to money laundering. A particularly important step was made with the Third Directive, after the September 11, 2001 terrorist attacks in the US, which stressed the problem of organized terrorism even more and made clear that tighter measures were needed at global level. Relevant advancements were made in January 2020 with the transposition by all member states of the Fifth AML Directive which aimed at expanding the extent of regulation also to virtual currency exchanges, estate agents, art dealers and more (LSEG, 2024) (European Union, 2018). The Sixth Directive came into effect in December 2020 not long after the previous one, and again, some major improvements were made. Specifically, they reached an harmonized definition for "predicate crime" ("cyber crime" and "environmental crime" had been included within the offences); they further expanded the regulatory scope also considering "aiding and abetting" as punishable criminal offences; they extended the criminal liability, so that companies could also be criminally liable for the illicit actions carried out by their employees; and lastly they also made punishments tougher by increasing the sentences for money laundering crimes.
The last Directive (EU) 2024/1640, known as AMLD VI, was adopted in April 2024 and represents a major overhaul of the EU's anti-money laundering (AML) framework. It builds upon and significantly strengthens previous directives (AMLD IV & V) through more centralized oversight, broader scope, and enhanced transparency. AMLD VI is the most ambitious and comprehensive EU AML directive to date. It establishes a new centralized authority (AMLA), mandates greater transparency and access to financial data, introduces tighter controls on crypto and cash, and broadens the regulatory scope to new sectors like football and luxury goods (see paragraph 3.6). It marks a major step toward a fully harmonized and enforceable EU-wide AML framework.

## 3.3) United Nations Convention Against Corruption (UNCAC)

Adopted in 2003 and entered into force in 2005, the UNCAC (also known as "Merida Convention") had the aim to promote preventing measures against corruption and the criminalization of acts like money laundering, bribery and embezzlement (UNODC, 2000). It is the only legally binding universal tool to contrast corruption. It requires member states to implement a set of anti-corruption measures and promotes both international cooperation and mutual legal assistance in the fight against illicit activities. The convention also plays a relevant role in driving forward both the 2030 Agenda and the Sustainable Development Goals (SDGs), by addressing the widespread problem of corruption (United Nations, 2021).

## 3.4) Bank Secrecy Act (BSA)

Overseas, the first law to combat money laundering was enacted as early as 1970 by the US Congress with the Bank Secrecy Act, officially known as the Currency and Foreign Transaction Reporting Act (IRS, 2025). Shortly after its passage, the BSA met the indignation of several groups who thought it was unconstitutional, claiming it was violating both the Fourth and Fifth Amendments. For these reasons it remained inactive until the '80s, when financial institutions eventually complied with BSA requirements.
Nowadays BSA mandates financial institutions to maintain records and submit different types of reports based on the problem encountered. It is a cornerstone regulation in the combat against money laundering and other fraudulent activities in the US.

## 3.5) Recent developments in Italian regulations (UIF)

On July 3, 2025, the Italian Financial Intelligence Unit (UIF) published a consultation paper updating the instructions regarding the reporting of suspicious transactions (SOS), with the aim of enhancing the effectiveness of the anti-money laundering and counter-terrorism financing system. The new text is intended to replace the current regulation dated May 4, 2011, in light of recent regulatory developments and international best practices. UIF seeks to improve the quality of reports by discouraging overly automatic or overly cautious reporting practices, which risk undermining the investigative value of suspicious transaction reports.
The instructions are divided into three main sections. The first part sets out the general principles and operational rules, emphasizing the importance of active cooperation by obligated entities. It clarifies that a report should not be triggered by numerical thresholds or automated criteria alone, but rather by a concrete, documented, and reasoned assessment of objective or subjective anomalies. The analysis process must be thorough, and in some cases, may include the temporary suspension of the suspicious transaction. Additionally, UIF encourages feedback mechanisms to strengthen the overall effectiveness of the reporting system. The second part of the document addresses the organizational and procedural obligations of reporting entities. Each organization or professional must

appoint a person responsible for SOS reporting, who must be independent, competent, and free from conflicts of interest. In smaller settings, such as individual practices, this responsibility may fall onto the professional directly. Entities are required to establish formal procedures for detecting and assessing suspicious activity, even when using IT tools or artificial intelligence algorithms. However, these tools must complement human analysis rather than replace it.

The third part focuses on the technical and operational aspects of submitting reports through the Infostat-UIF portal. It provides guidelines on how to register, compile, submit, amend, or cancel reports, with the goal of streamlining the process and ensuring consistent information flows.

The document has been opened to public consultation for a period of 60 days, ending on September 3, 2025. UIF invites all relevant stakeholders, including court-appointed administrators, to submit comments and suggestions. Particular attention is given to those operating in high-risk contexts, such as the management of seized or confiscated assets.

Overall, UIF's proposed reform marks a decisive step toward a more selective, professional, and effective reporting system. The goal is not to increase the number of reports, but to enhance their quality, favoring thoughtful analysis over automatic or overly cautious reporting. This change requires the active commitment of all obligated entities, who are called upon to exercise sound and responsible judgment in managing money laundering and terrorism financing risks (UIF, 2025).

## 3.6) Toward centralized AML supervision in Europe

In 2024, the European Union introduced a major reform of its anti-money laundering framework with the approval of a new legislative package, including two Regulations and one Directive. Most notably, Regulation (EU) 2024/1620 established the European Anti-Money Laundering Authority (AMLA), headquartered in Frankfurt. AMLA will directly supervise selected high-risk financial institutions and coordinate national AML/CFT authorities, with the goal of harmonizing and strengthening supervisory practices across the EU (European Union, 2024 [a]). Alongside this institutional innovation, Regulation (EU) 2024/1624 and Directive (EU) 2024/1640 lay out new rules on risk assessment, customer due diligence, and cross-border cooperation. This transition to centralized supervision marks a clear shift toward greater efficiency, consistency, and technological adoption in the fight against financial crime.

The approach proposed in this study, combining Machine Learning techniques with expert-driven heuristic rules, aligns with the European regulatory direction, which increasingly promotes data-driven supervision and enhanced detection capabilities supported by innovation and harmonized methodologies (European Union, 2024 [b]) (European Union, 2024 [c]). Furthermore, the study is particularly timely considering the recent AML regulatory developments proposed by UIF whose documents are discussed by professionals in this field.

## 4) Description of the methodologies

In this section, we analyze in detail the three predictive models used in the study, namely: Decision tree, Logistic and Logic Learning Machine. Each model takes a different approach to classification, with distinct characteristics that influence both performance and interpretability of the results. The underlying idea is to check how the Logic Learning Machine performs and then compare its results with the ones obtained through the other two standard classification techniques. Therefore, we first describe the working principles of the models considered, analyzing their key features, in such a way as to provide a clear understanding of each algorithm. Beyond describing the models, we also analyze the evaluation metrics used to assess their effectiveness and determine the best configuration.

## 4.1) Traditional Machine Learning methodologies

Traditional Machine Learning methodologies, such as Classification and Regression Trees (CART) and logistic regression, have long been foundational in data analysis, providing powerful tools for extracting patterns and making predictions.

CART models represent a method for approximating data through a stepwise function, obtained by iteratively subdividing the space of observations, based on specific thresholds applied to the model variables. Each subdivision aims to identify subsets of data with values of the target variable that are as homogeneous as possible. In classification problems, observations are assigned to the most representative class within the group to which they belong. This methodology is often visualized as a decision tree, since the separation rules can be organized in a hierarchical structure. Each subdivision of the dataset can be represented as a node in a binary tree, where the first node, known as the "root", serves as the starting point, while the terminal nodes, called "leaves", identify the final subsets of data. To train a decision tree, a dataset is needed where the target variable is known. The algorithm builds the model progressively, applying binary splits to the data, based on optimal cutoffs. In each phase, a threshold is determined for each variable which allows to obtain a subdivision such as to reduce the variance within each subset and, at the same time, increase its difference compared to the other groups. This procedure is repeated iteratively on each new subset, determining, each time, the optimal threshold for the variable considered. The process continues until a stopping criterion is reached, which may be, for example, the impossibility of further improving the separation of the data, the presence of only one element in a subset or the creation of a group composed exclusively of observations belonging to the same class as the target variable.

Once training is complete, the resulting tree can become very complex and detailed, with a large number of splits that make it difficult to interpret and increase the risk of over-fitting. To avoid this problem, a simplification technique known as "pruning" is applied. This process is guided by a loss function which allows to reduce the complexity of the tree by removing less significant branches. Pruning is performed by progressively eliminating branches that contribute the least to the overall performance, aiming to achieve a good balance between model simplicity and predictive accuracy (Lewis, 2000).

Logistic is a supervised algorithm used for binary classification tasks. It is widely employed in areas such as fraud detection, medical diagnosis and engineering (e.g. for predicting the probability of failure of a specific process). It is used to model the probability that a given instance belongs to one of two classes of the binary categorial dependent variable. This method uses the logistic function, which is able to convert real values into an interval between 0 and 1: this ensures that the predicted probabilities are in this range

(Ohno-Machado et al., 2002). The logistic function is of the form: $p(x) = \frac{1}{1+\exp(-(x-\mu)/s)}$ where $\mu$ is a location parameter (i.e. the midpoint of the curve, where $p(\mu) = 0.5$) and $s$ is a scalar parameter.

## 4.2) Logic Learning Machine (LLM)

The Logic Learning Machine (LLM) is a rule-based method alternative to decision trees. In plain words, the LLM transforms the data into a Boolean domain where some Boolean functions (namely one for each output value) are reconstructed starting from a portion of their truth table with a method that is described in the paper of Muselli and Ferrari (Muselli et al., 2011). The method creates a set of intelligible rules through Boolean function synthesis following 4 steps. These steps are:

1. Discretization
2. Latticization or Binarization
3. Positive Boolean function
4. Rule generation

In a classification problem, $d$-dimensional examples $x \in X \subset \mathfrak{R}^d$ are to be assigned to one of $q$ possible classes, labeled by the values of a categorical output $y$. Starting from a training set $S$ including $n$ pairs $(x_i, y_i), i = 1, \dots, n$, deriving from previous observation, techniques for solving classification problems have the aim of generating a model $g(x)$, called classifier, that provides the correct answer $y = g(x)$, for most input patterns $x$. In order to analyze the process, a bi-class toy problem is used, whose training set is shown in Table 1. In this example $O_0$ represents a normal transaction, whilst $O_1$ represents a fraudulent transaction.

| $X_1$ | 700 | 1100 | 2200 | 1400 | 2300 | 800 | 1200 | 2100 | 2600 | 2400 |
|---|---|---|---|---|---|---|---|---|---|---|
| $X_2$ | Cheque | Cheque | Cash withdrawal | Cheque | Credit card | Cash withdrawal | Credit card | Cheque | Cheque | ACH |
| $Y$ | $O_0$ | $O_0$ | $O_0$ | $O_0$ | $O_0$ | $O_1$ | $O_1$ | $O_1$ | $O_1$ | $O_1$ |

*Table 1: Toy example for describing the LLM working principle*

## 4.2.1) Discretization

In this step, each continuous variable domain is converted into a discrete domain by a mapping. $\psi_j X: X_j \rightarrow I_M$ where $X_j$ is the domain of the $j$-th variable and $I_M = 1, \dots, M$ is the set of positive integers up to $M$. The mapping must preserve the ordering of the data. If $x_{ij} \leq x_{kj}$ then $\psi_j(x_i) \leq \psi_j(x_k)$, $\forall j = 1, \dots, d$. One way to describe $\psi_j$ is that it consists of a vector $\boldsymbol{\gamma}_j = (\gamma_{j1}, \dots, \gamma_{jm}, \dots, \gamma_{M_j-1})$ such that:

$$\psi_j(\boldsymbol{x}_i) = \begin{cases} 1, & x_{ij} \leq \gamma_{j1} \\ m, & \gamma_{jm-1} < x_{ij} \leq \gamma_{jm} \\ M_j, & x_{ij} > \gamma_{jM_j-1} \end{cases} \quad \text{(1)}$$

There are several strategies for discretization and the simplest one is creating $M_j$ interval having the same length. Let $\boldsymbol{\rho}_j$ be the vector of all the $\alpha_j$ values for input variable $j$ in ascending order $(p_{jl} < p_{jl+1} \ \forall \ l = 1, \dots, \alpha_j)$, then the cutoff $\gamma_{jm}$ is given by:

$$\gamma_{jm} = p_{j1} + \frac{p_{j\alpha_j} - p_{j1}}{M_j} m \quad \text{(2)}$$

This method is referred to as Equal Width discretization. Another approach defines one interval for each value.

## 4.2.2) Binarization

In this step, each discretized domain is transformed into a binary domain through a mapping $\varphi_j: I_{M_j} \rightarrow \{0,1\}^{M_j}$, where $I_{M_j}$ is the domain of the $j$-th variable and $\{0,1\}^{M_j}$ is a string having a bit for each possible value in $I_{M_j}$. The mapping must maintain the ordering of data: $u < v$ if and only if $\varphi_j(u) < \varphi_j(v)$ where the standard ordering between $\boldsymbol{z}$ and $\boldsymbol{w} \in \{0,1\}^{M_j}$ is defined as follows:

$$\boldsymbol{z} < \boldsymbol{w} \text{ if and only if } \begin{cases} \exists \ i \quad \text{such that } z_i < w_i \\ \forall l \ < i \quad z_l \leq w_l \end{cases}$$

$$\text{(3)}$$

$$\boldsymbol{z} \leq \boldsymbol{w} \text{ if and only if } z_i \leq w_i \ \forall \ i = 1, \dots, M_j$$

if the relation in the equation above holds, then it is said that $\boldsymbol{z}$ covers $\boldsymbol{w}$.

A suitable choice for $\varphi_j$ is the inverse only-one coding, that for each $k \in I_{M_j}$ creates a string $\boldsymbol{h} \in \{0,1\}^{M_j}$ having all bits equal to 1 except the $k$-th bit which is set to 0. For example, let $x_{ij} = 3$ with domain $I_5$, then $\varphi_j(\boldsymbol{x}_i) = 11011$. In this way $\varphi_j(\boldsymbol{x}_i) = \boldsymbol{z}_i$ where $\boldsymbol{z}_i$ is obtained by concatenating $\varphi_j(\boldsymbol{x}_i)$ for $j = 1, \dots, d$. As a result, the new training set is $S' = \{(\boldsymbol{z}_i, y_i)\}_{i=1}^N$, with $\boldsymbol{z}_i \in \{0,1\}^B$ where $B = \sum_{j=1}^d M_j$. The training set obtained by applying discretization with the single cutoff 1500 for the variable $X_1$ and subsequent binarization for the toy problem is shown in Table 2.

| Z | Y |
|---|---|
| 01 0111 | Normal |
| 01 0111 | Normal |
| 10 1101 | Normal |
| 01 0111 | Normal |
| 10 1110 | Normal |

| Z | Y |
|---|---|
| 01 1101 | Fraud |
| 01 1110 | Fraud |
| 10 0111 | Fraud |
| 10 0111 | Fraud |
| 10 1011 | Fraud |

*Table 2: Toy example after binarization.*

## 4.2.3) Synthesis of the Boolean function

The training set $S'$, obtained after binarization, can be divided into two different subsets according to the output class: $T$ is the set containing $(\boldsymbol{z}_i, y_i)$ with $y_i = O_1$ whereas $F$ is the set containing the example for which $y_i = O_0$. $T$ and $F$ can be viewed as a portion of the truth table of a Boolean function $f$ that must be reconstructed. Before proceeding with the method description, it is useful to give some definitions and notations.

- Each Boolean function can be written with operators AND, OR, and NOT that constitute the Boolean algebra; if NOT is not considered then a simpler structure, called Boolean lattice, is obtained. From now on, only the Boolean lattice is considered. It can be drawn by positioning $\boldsymbol{z}$ over $\boldsymbol{w}$ if $\boldsymbol{z} > \boldsymbol{w}$ and by linking all the couples $\boldsymbol{z}$, $\boldsymbol{w}$ for which an $\boldsymbol{a}$ does not exist such that $\boldsymbol{w} < \boldsymbol{a} < \boldsymbol{z}$. An example for $\{0,1\}^3$ is shown in Figure 1.b.

- The sum (OR) and product (AND) of $\eta$ terms can be denoted as follows:

$$\bigvee_{j=1}^{\eta} z_j = z_1 + z_2 + \cdots + z_\eta \quad (4)$$
$$\bigwedge_{j=1}^{\eta} z_j = z_1 \cdot z_2 \cdot \dots \cdot z_\eta = z_1 z_2 \dots z_\eta$$

- A logical product is called an implicant of a function $f$ if the following relation holds: $\bigwedge_{j=1}^{\eta} z_j \leq f$, where each element $z_j$ is called literal. The product is called prime implicant if the relation no longer holds when a literal is removed from the implicant.

- The ordering in a Boolean lattice is defined by the equations above; according to this ordering, a Boolean function $f: \{0,1\}^B \rightarrow \{0,1\}$ is called positive if $\boldsymbol{z} \leq \boldsymbol{w}$ implies $f(\boldsymbol{z}) \leq f(\boldsymbol{w})$ for each $\boldsymbol{z}, \boldsymbol{w} \in \{0,1\}^B$.

- A subset $A \subset I_B$ such that for each element $\boldsymbol{z}, \boldsymbol{w} \in A$, an ordering cannot be established (neither $\boldsymbol{z} < \boldsymbol{w}$, nor $\boldsymbol{w} < \boldsymbol{z}$), is called antichain.

- Given $\boldsymbol{a} \in \{0,1\}^B$, then the set $L(\boldsymbol{a}) = \{\boldsymbol{z} \in \{0,1\}^B \mid \boldsymbol{z} \leq \boldsymbol{a}\}$ is called lower shadow of $\boldsymbol{a}$, whereas the set $U(\boldsymbol{a}) = \{\boldsymbol{z} \in \{0,1\}^B \mid \boldsymbol{z} \geq \boldsymbol{a}\}$ is called an upper shadow of $\boldsymbol{a}$. The lower and upper shadows for $101 \in \{0,1\}^3$ are shown in Figures 1.a and 1.c.

- Given the subset $T, F \in \{0,1\}^B$, then $T$ is lower separated from $F$, if there is no element $\boldsymbol{z} \in T$ belonging to the lower shadow of some element of $F$.

- Given the binary string $\boldsymbol{a}$, if there is a $\boldsymbol{z} \in T$ such that $\boldsymbol{a} \leq \boldsymbol{z}$, there is not a $\boldsymbol{w} \in F$ such that $\boldsymbol{a} \leq \boldsymbol{w}$, and for each $\boldsymbol{b} < \boldsymbol{a}$, there is $\boldsymbol{w} \in F$ such that $\boldsymbol{b} \leq \boldsymbol{w}$, then $\boldsymbol{a}$ is called bottom point for the pair $(T, F)$.

- Every positive Boolean function can be written in its unique, not redundant Positive Disjunctive Normal Form (PDNF) as the sum of its prime implicants: $f(\boldsymbol{z}) = \bigvee_{\boldsymbol{a} \in A} \bigwedge_{j \in P(\boldsymbol{a})} z_j$, where $P(\boldsymbol{a})$ is the subset $I_B$ containing each $i$ such that $a_i = 1$; $A$ is an antichain of $\{0,1\}^B$ and each $\boldsymbol{a}$ is called the minimum true point. For example, the not redundant PDNF $f(\boldsymbol{z}) = z_1 z_3 + z_4$ is obtained from antichain $A = \{1010, 0001\}$.

From these definitions, it follows that a method for finding $f$ must retrieve the set of minimum true points to be used from $T$ and $F$, in order to represent $f$ in its irredundant PDNF and it follows that the set of all bottom points for $(T, F)$ is an antichain, which elements are candidate minimum true points.

The algorithm employed by LLM to produce implicants is called Shadow Clustering (Muselli et al., 2011). It generates implicants for $f$ through the analysis of the Boolean lattice $\{0,1\}^B$. The algorithm selects a node in the diagram and generates bottom points $(T, F)$

by descending the diagram: moving down from a node to another node is equivalent to changing a component from 1 to 0 and a bottom point is added to A when any further move down leads to a node belonging to the lower shadow of some $w \in F$.

In particular, the starting node is chosen between the $z \in T \subset \{0,1\}^B$ that do not cover any point $a \in A$ such that $a \leq z$ (in other words the algorithm ends when each element in $T$ covers at least one element in $A$). Once $A$ has been found, it is possible that it contains redundant elements and consequently, it must be simplified in order to find $A^*$, from which the PDNF of the positive Boolean function $f$ can be derived.
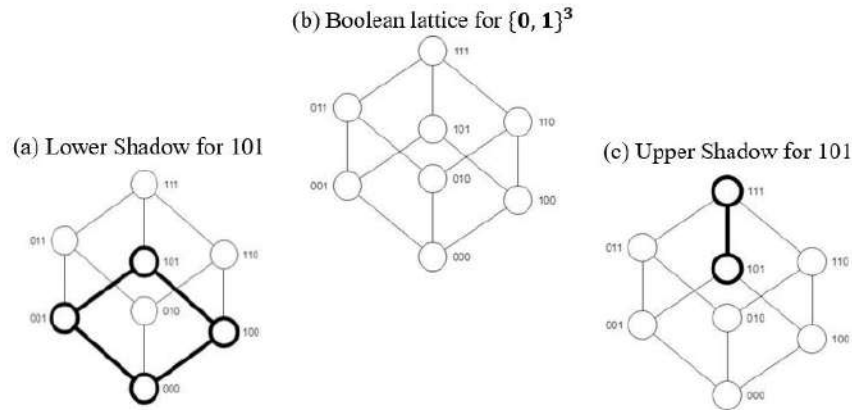


*Figure 1: Boolean Lattice diagram. Source: Muselli et al., 2011*

Different versions of Shadow Clustering exist depending on the choice of the element to be switched from 1 to 0 at each step of the diagram descent. For example, Maximum-covering Shadow Clustering (MSC) at each step changes the $i$-th element that maximizes the associated potential covering, defined as the number of elements $z \in T$ for which $z_i = 0$.

As concerns the selection of $A^* \subset A$, a possible choice is to subsequently add to $A^*$ the element of $A$ that covers the highest number of points in $T$ that are not covered by any other element of $A^*$. The application of the Shadow Clustering algorithm (Algorithm 1) to the dataset after binarization shown in Table 2 produces the implicants:

$$100011 \text{ which corresponds to } z_1 \wedge z_5 \wedge z_6$$
$$011100 \text{ which corresponds to } z_2 \wedge z_3 \wedge z_4$$

Then, the PDNF of the resulting Boolean function is the following: $f(z) = (z_1 \wedge z_5 \wedge z_6) \vee (z_2 \wedge z_3 \wedge z_4)$

---

**Algorithm 1** Shadow Clustering algorithm (bottom-up)

---

**Data**: $P(x)$
$I = P(x)$;
$A = \emptyset$;
**while** $I \neq \emptyset$ **do**
    choose $i \in I$ and remove it from $I$
    if there is $y \in F$ such that $p(I \cup A) \leq y$ then add $i$ to $A$
**end**
**Return** $p(A)$.

---

## 4.2.4) Rule generation

In the last step, each implicant of the positive Boolean function $f$ is transformed into an intelligible rule, where, as said before, a function is generated for each output value, and then the consequent of the rules only depends on $f$. The transformation considers the coding applied during binarization. In particular, $z$ was obtained by concatenating the results of the mapping $\varphi_j(x)$ for each $j = 1, \dots, d$ and consequently it can be split into substring $h_j$ for each attribute, whose bit $z_i \in h_j$ corresponds to a nominal value if $X_j$ is nominal, whereas it corresponds to an interval if $X_j$ is ordered.

For each implicant, a rule in IF − THEN form is generated by adding a condition for each attribute $X_j$ as follows:

- If $z_i = 0$ for each $z_i \in h_j$, then no condition relative to $X_j$ is added to the rule;

- If $X_j$ is nominal, then a condition $X_j \in V$ is added to the rule, where $V$ is the set of values associated with each $z_i \in h_j$ such that $z_i = 0$;

- If $X_j$ is ordered, then a condition $X_j \in V$ is added to the rule, where $V$ is the union of the intervals associated with each $z_i \in h_j$ such that $z_i = 0$.

For the implicant 100011 obtained in the previous step, $h_1 = 10$ leads to the condition $X_1 \in (1700, inf)$ or $X_1 > 1700$ and $h_2 = 0011$ leads to the condition $X_2 \in \{Check, ACH\}$. Then the rule relative to 100011 is: IF $X_1 > 1700$ AND $X_2 \in \{Check, ACH\}$ THEN $Y = O_1$

For the implicant 011100 obtained in the step described previously, $h_1 = 01$ leads to the condition $X_1 \in (-inf, 1700]$ or $X_1 \leq 1700$ and $h_2 = 1100$ leads to the condition $X_2 \in \{Cash\ withdrawal, Credit\ card\}$. Then the rule relative to 011100 is: IF $X_1 \leq 1700$ AND $X_2 \in \{Cash\ withdrawal, Credit\ card\}$ THEN $Y = O_1$

If $X_j$ is ordered, conventionally the upper bound of the interval, if finite, is always included in the condition, whereas the lower bound is excluded. In order to generate the rule for the other class, it is sufficient to label $O_0$ with 1 and $O_1$ with 0. In the case of the multiclass problem, it is sufficient to decompose the problem into several bi-class problems for each of the sub-problems the target class is labelled with 1, and all the remaining with 0.

## 4.2.5) Rule quality and class prediction

The process described in the previous subsection implies that each element $x_i$ of the training set only satisfies rules associated with the output class of $x_i$, but since data are affected by noise, usually it is preferable to admit some errors in order for the model to be able to generalize. In order to permit a fraction of error, the descent of the diagram does not stop when a further move down leads to the lower shadow of some $w \in F$, but still allowing it to go on until a further move leads to a node belonging to the lower shadow of a percentage element $w \in F$ greater than a regularization parameter $\varepsilon_{max}$. Then, it is usual that an element of the training set covers the rule of different classes. When it happens, the output class is established according to the relevance of the rules satisfied by it.

In order to present relevance, the following quantities relative to a rule $r$ in the IF $< premise >$ THEN $< consequence >$ form are introduced:

- $TP(r)$ is the number of training set examples that satisfy both the premise and the consequence of the rule $r$;
- $FP(r)$ is the number of training set examples that satisfy the premise but do not satisfy the consequence of the rule $r$;
- $TN(r)$ is the number of training set examples that do not satisfy either the premise or the consequence of the rule $r$;
- $FN(r)$ is the number of training set examples that do not satisfy the premise and satisfy the consequence of the rule $r$.

Please note that an example $x_i$ satisfies the premise of the rule $r$ if it satisfies all its premise conditions, whereas $x_i$ does not satisfy the premise of the rule $r$ if it does not satisfy at least one among its premise conditions. Combining these quantities, it is possible to compute quality measures for a rule $r$:

Covering: $C(r) = \frac{TP(r)}{TP(r)+FN(r)}$ (5)

Error: $E(r) = \frac{FP(r)}{TN(r)+FP(r)}$ (6)

It is evident that the greater the covering, the more relevant the rule is; on the other hand, the smaller the error, the less relevant the rule is. Then, the relevance of a rule $r$ is obtained by combining $C(r)$ and $E(r)$: $R(r) = C(r)(1 - E(r))$.

Once the relevance of the rule is defined, it is possible to use it to compute a score $S(x_i, o)$ for each class $o$ that measures how likely it is that $y_i = o$:

$$S(x_i, o) = \sum_{r \in \mathcal{R}_o^i} R(r) \text{ (7)}$$

with $\mathcal{R}_o^i = \{r \mid r \in \mathcal{R}, r \leq x_i, O(r) = o\}$, where $\mathcal{R}$ is the complete ruleset, $r \leq x_i$ denotes that $x_i$ satisfies the premise of the rule. On the other hand, to obtain a measure of relevance $R(c)$ for a condition $c$ included in the premise part of a rule $r$, the rule $r'$ can be considered, obtained by removing that condition from $r$. Since the premise part of $r'$ is less stringent, we obtain that $E(r') \geq E(r)$ so that the quantity $R(c) = (E(r') - E(r))C(r)$ can be used as a measure of relevance for the condition $c$ of interest. $O(r) = o$ denotes the consequence of $r$ predict class $o$. Then $\mathcal{R}_o^i$ is the set of rules satisfied by $x_i$ that predict class $o$. From the scores of each output class, it is possible to define the probability that $y_i = o$:

$$P(o \mid x_i) = \frac{S(x_i, o)}{\sum_{k \in O} S(x_i, k)} \text{ (8)}$$

Then the selected output is the one that maximizes the output probability: $\tilde{y}_i = \max_o P(o \mid x_i)$

## 4.2.6) Feature ranking

For every ordered variable $x_j \in Z$, let us denote with $M_j - 1$ the collection of all the thresholds $\gamma_{jl}$ involved in the conditions of rules $r_k$; through these thresholds the domain of the component $x_j$ is subdivided into $M_j$ adjacent intervals $[-\infty, \gamma_{j1}], (\gamma_{j1}, \gamma_{j2}], \ldots, (\gamma_{j,l-1}, \gamma_{jl}], \ldots, (\gamma_{j M_j-1}, +\infty]$. Let us denote with $J_{j1}, J_{j2}, \ldots, J_{M_j}$ these intervals, so that $J_{j1} = [-\infty, \gamma_{j1}]$, $J_{j2} = (\gamma_{j1}, \gamma_{j2}]$, etc.

Now, if a rule $r_k \in \mathcal{R}_o = \{r \mid r \in \mathcal{R}, O(r) = o\}$ for the output class $o$ (i.e. whose consequence part is $y = o$) includes a condition $c_{kl}$, with relevance $R(c_{kl})$, involving the ordered component $x_j$, the points of $m_{kl}$ of the $M_j$ adjacent intervals verify that condition. For instance, if the condition $c_{kl}$ is $x_j \leq \gamma_{j3}$, the points of the $m_{kl} = 3$ intervals $J_{j1}$, $J_{j2}$ and $J_{j3}$ satisfy $c_{kl}$. It is then possible to retrieve a measure of relevance $R_k^o(J_{ji})$ for each interval $J_{ji}$, with respect to the output class $o$, by looking at the quantities $R(c_{kl})$ of the conditions $c_{kl}$, that are included in rules $r_k$, that involve the component $x_j$, and are verified by points of $J_{ji}$. In particular, if a condition $c_{kl}$ involving $x_j$ is satisfied by $m_{ki}$ of the $M_j$ adjacent intervals, the relevance quantity that can be attributed to each of these intervals is $R_k^o(J_{ji}) = R(c_{kl})/m_{kl}$.

By collecting all the relevancies derived from all the rules $r_k \in \mathcal{R}_o$ including a condition $c_{kl}$ on the component $x_j$, we can obtain the measure of relevance $R^h(J_{ji})$ of the interval $J_{ji}$ with respect to the output class $o$:

$$R^o(J_{ji}) = 1 - \prod_{r_k \in \mathcal{R}_o} \left(1 - R_k^o(J_{ji})\right) \quad (9)$$

Starting from Eq. (9) a measure of relevance $R^o(x_j)$ for the component $x_j$ (with respect to $o$) can be derived by considering the variation of $R^o(J_{ji})$ over the $M_j$ adjacent intervals $J_{j1}$, $J_{j2}$, ..., $J_{M_j}$. In fact, if $R^o(J_{ji})$ does not change so much in these intervals, then different thresholds are essentially used to determine parts of the input domain where the behavior of the model $g(x)$ is similar. This means that the variable $x_j$ has little discriminant power among different classes, but it characterizes the input domain with respect to $g(x)$ for the output class $o$.

A possible way of measuring the variation of a quantity is to consider its standard deviation $\sigma$; therefore, we have:

$$R^o(x_j) = M_j\, \sigma_j \left(R^o(J_{ji})\right) \quad (10)$$

where $\sigma_j$ stands for the standard deviation over the $M_j$ intervals $J_{j1}$, $J_{j2}$, ..., $J_{M_j}$.

A sign for $R^o(x_j)$, which indicates if the variable $x_j$ is directly (if the sign is positive) or inversely (if the sign is negative) correlated with the output class $o$, can also be retrieved by looking where higher values of $R^o(J_{ji})$ are located. In particular, if higher values of $R^o(J_{ji})$ occur at higher (resp. lower) $i$ then the variable $x_j$ is directly (resp. inversely) correlated with the output class $o$.

Hence, a procedure for deriving the sign of $R^o(x_j)$ consists in subdividing the product of (1) in two parts: the first one, denoted with $R^{o-}(J_{ji})$, contains terms $R_k(J_{ji})$ originated by conditions $c_{kl}$ of the form $x_j \leq \gamma_{ji}$, whereas $R^{o+}(J_{ji})$ includes terms $R_k(J_{ji})$ derived by conditions $c_{kl}$ of the kind $x_j > \gamma_{ji}$. As for conditions $c_{kl}$ of the form $\gamma_{ji_1} < x_j \leq \gamma_{ji_2}$ terms $R_k(J_{ji})$ for $i \leq (i_1+i_2)/2$ (resp. $i > (i_1+i_2)/2$) are inserted into $R^{o-}(J_{ji})$ (resp. $R^{o+}(J_{ji})$). With these definitions, the sign of $R^o(x_j)$ becomes negative if $R^{o-}(J_{ji}) < R^{o+}(J_{ji})$ and positive in the opposite case.

If the variable $x_j$ is nominal, then equation (1) can still be used to determine measures of relevance $R^o(J_{ji})$ if $G_j = \{v_{j1}, v_{j2}, ...\}$ is the collection of the possible values assumed by $x_j$ and $J_{ji} = \{v_{ji}\}$, for $i = 1, 2, ..., |G_j|$. In this case equation (9) becomes:

$$R^o(x_j) = |G_j| \sigma_j \left(R^o(J_{ji})\right) \quad (11)$$

a sign for $R^o(x_j)$ cannot be determined and is therefore always considered as positive.

If the (absolute) maximum over the $q$ output classes of the quantities $R^o(x_j)$ is greater than 1, then all the relevancies $R^o(x_j)$ are normalized to this maximum so that their values lie in the range $[0,1]$. By averaging the quantities $R^o(J_{ji})$ and $R^o(x_j)$, for $o = 1, ..., q$, we can obtain absolute measures of relevance $R(J_{ji})$ and $R(x_j)$ (independent of the output class $o$) for $J_{ji}$ and for the variable $x_j$:

$$R(J_{ji}) = \frac{1}{q}\sum_{o=1}^{q} R^o(J_{ji}) \quad , \quad R(x_j) = \frac{1}{q}\sum_{o=1}^{q} R^o(x_j) \quad (12)$$

In short, as regards Logic Learning Machine models, feature importance can be analyzed based on the generated rules and their frequency and predictive strength. In fact, the Logic Learning Machine does not use coefficients such as the Logistic task or the number of the splits like the Decision tree task, but instead, its feature importance is based on the relevance of the features in the ruleset extracted from the model. By using feature ranking applied to LLM, it is possible to inspect the presence and weight of attributes in the final ruleset. This task provides an analysis of the feature importance by counting how many times a feature appears in the model rules and generates a ranking of features, showing which ones had the greatest impact. The more a feature appears in important rules, the more impact it has on model decisions.

## 4.3) Evaluation metrics

As previously stated, the objective of this study was to confirm whether the combination of Machine Learning algorithms with heuristic rules could lead to an improvement of the results in the detection of fraudulent transactions. Therefore, our aim was to analyze and find sensible heuristic rules that could somehow bring an improvement both in terms of precision and explicability of the results. In writing and selecting the rules, we therefore tried to combine the information deriving from regulations and from the SAML-D paper, to obtain a general picture of the topic and to summarize this information within our personal ruleset.

In order to evaluate the performance of the heuristic rules, we employed the "covering" and "error" statistics previously analyzed (see Eq. 5 and 6). The covering describes the percentage of samples that are covered by that rule, in a class, compared to the total samples in that class. We want this value to be the highest possible, i.e. $\approx 1$. The error, on the contrary, measures the percentage of errors within the covering of the rule, i.e. how often the rule is wrong within the covering. In this case we aim to obtain a small error, i.e. $\approx 0$. These statistics were computed to later calculate further metrics that allowed us to select only the best performing rules and group them together.

In this phase, the preliminary step was to compute specific metrics in order to filter and pick only the best rules for each of them. To this aim, we used the error and covering statistics we previously calculated. Specifically, we computed the following metrics:

- **Error/Covering**: this metric simply performs the ratio between the two core statistics. It indicates the proportion of error compared to how much the rule is applicable in the dataset. Since we wanted our rules to maximize the covering and minimize the error, we consequently selected only those rules who scored the lowest results for this metric;

- **Score**: this metric is useful for balancing precision and generalization. Specifically, it "discounts" the covering for the error. The formula to describe it is:

$$Score = Covering \cdot (1 - Error) \quad (13)$$

- **Score (1)**: this is an alternative version of the previous score metric. In this case we reduce the weight of the covering based on the ratio between error and covering. The formula to describe it is:

$$Score_1 = Covering \cdot \left(1 - {Error}/{Covering}\right) \quad (14)$$

The aim at this point was to rank the rules according to the metrics just described, to select only the best performing ones for each of them. Consequently, we filtered the first 10 rules which scored the best and selected them to later complement the classification models in the merging phase between heuristics and Machine Learning.

As an example of one of the best-performing rules identified, we present the "*AMLCheckUAE*" rule. This rule was selected based on its high covering and low error, reflecting its effective application in detecting potentially fraudulent transactions. The rule flags transactions as fraudulent if either the sender's or receiver's bank location is in the United Arab Emirates (UAE), and if the transaction amount exceeds the defined threshold for specific payment types. The logic behind this rule stems from the known risks associated with high-value transactions in particular locations and payment methods. By applying this rule, we are able to identify potentially suspicious transactions and flag them as fraudulent.

In general, the results seem to show that the best rules obtained are quite specific, in the sense that they capture very few fraud cases in the dataset. This is not necessarily a bad thing but this lack in intercepting fraudulent instances may indicate a scarce relevance in terms of improvement of results. In other words, the most prominent heuristic rules we defined are generally very small, that is, they have a negligible weight compared to what classification algorithms can achieve (they are often precise but have little generality).

To compare the employed classification models' performance, we adopted several evaluation metrics. Each of them provides specific information on the prediction quality, allowing for an extensive analysis of the results. This was done to discriminate against the model and the specific parameterization which performed best among all the ones we implemented.

- **AUC**: The AUC (Area Under the Curve) measures the ability of a model to distinguish between two target classes, in this case between being a laundering transaction or a normal one. It is computed exploiting the Receiver Operating Characteristic (ROC) curve, which represents, for different thresholds, the tradeoff between the True Positive Rate (TPR), also known as sensitivity, and the False Positive Rate (FPR). In this study we have implemented it by applying the roc function to the score of the model predictions. The results in terms of AUC vary in the range [0.5, 1], where the left bound indicates a non-discriminating model, while the right one denotes a perfectly discriminating model.

- **Precision**: The Confusion matrix shows the number and percentage value of correctly and incorrectly classified observations. In this context, the precision statistic measures the proportion of correct positive predictions with respect to the total number of positive predictions made by the model. In mathematical terms, this is defined as:

$$Precision = \frac{TP}{TP+FP} \quad (15)$$

High values of this metric indicate that the model is effective at reducing false positives, meaning it rarely misclassifies normal transactions as fraudulent. Optimizing precision is crucial in fraud-detection analysis, mainly because false alarms generate a cost for the agency that controls suspicious cases. Therefore, the main goal is to reduce the number of false positives and improve the overall effectiveness of the model.

- **Youden J statistic**: This statistic measures a classification model's ability to discern between two classes. In mathematical terms, it is described as:

$$J = TPR + TNR - 1 \quad (16)$$

It combines the model sensibility, captured by the true positive rate, with its specificity, indicated by the true negative rate, into a synthetic index. Youden's J statistic has proven to be a very useful index for two main reasons. First of all because it is not affected

by class imbalance which, as we know, is a problem affecting SAML datasets, and second because it is only valid for binary classification problems. Its value ranges between [-1, 1], where J=0 indicates that the model has no discriminating ability (it has the same behavior as the random case); J=1 denotes a perfect classifying model which commits no errors; whilst negative values signal rare pathological cases where the model performs worse than the random case.

## 4.4) SAML Dataset at a glance

To address the analyzed problem, the dataset used was the Synthetic AML Dataset (SAML-D), available on Kaggle (Oztas et. al, 2023). SAML-D incorporates 12 features and 28 typologies of transactions (see Appendix B), split between 11 normal and 17 suspicious, making it one of the most comprehensive synthetic AML datasets available. These typologies have been selected based on existing datasets, academic literature, and interviews with AML specialists.

For the construction of SAML dataset certain rules and filters were implemented. In particular, the generation process of both "normal" and "suspicious" transactions involved two methods: the agent-based approach and the typology-based approach. This implies that the dataset includes prior assumptions about what constitutes normal and suspicious behaviors. For instance, fraudulent transactions are generated exploiting specific typologies like "structuring" or "deposit and send" which are characterized by specific patterns.

As a result, given the artificially encoded structure of these typologies, Machine Learning models trained on this dataset could partially learn to recognize these specific patterns rather than truly uncover novel laundering strategies. Consequently, this could limit the model's ability to generalize when applied to real-world data, which could feature new laundering behaviors not covered by the typologies simulated in the generation process (Oztas et. al, 2023).

In order to add complexity and realism to the data, observations include innovative features such as the geographic location of accounts, which also contains high-risk countries in the AML field, and high-risk payment types. In this sense, complexity and realism are also achieved by making fraudulent accounts carry out a wide range of money laundering types in addition to normal transactions (Oztas et al., 2023). Lastly, since the dataset was created with a focus on the United Kingdom, the prevalence of its observations, 99.72% if considering both inwards and outwards transactions, is therefore located in the UK.

The dataset comprises 9,504,852 transactions, of which 0.1039% are suspicious, thus showing a great class imbalance. This generally poses serious concerns in classification problems as the Machine Learning model implemented could develop bias towards the majority class, in this case non-fraud transactions. In other words, this can lead to a model that poorly learns the minority class, the one we are interested in, because it has few examples in the dataset.

When working with heavily unbalanced datasets, the absence of corrective actions can lead to sub-optimal results or misinterpretation of them. In particular, high false negative rates are likely to be obtained, as models tend to favor the majority class, partially or totally ignoring the minority one. This can make some common metrics, such as accuracy, unrepresentative of the model's real predictive capacity. Therefore, in these cases, it is essential to be aware of this issue and adopt appropriate metrics that take into account the imbalance. Moreover, it is also crucial to consider the impact of such an imbalance when comparing different models or between different parameterizations to avoid misleading conclusions based on poorly informed indicators.

| | Time | Date | Sender_account | Receiver_ac... | Amount | Payment_currency | Received_... | Sender_bank_l... | Receive... | Payment_type | Is_laundering | Laundering_type |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 10:35:19.0... | 2022-10-07 | 8724731955 | 2769355426 | 1459.150 | UK pounds | UK pounds | UK | UK | Cash Deposit | False | Normal_Cash_Deposits |
| 2 | 10:35:20.0... | 2022-10-07 | 1491989064 | 8401255335 | 6019.640 | UK pounds | Dirham | UK | UAE | Cross-border | False | Normal_Fan_Out |
| 3 | 10:35:20.0... | 2022-10-07 | 287305149 | 4404767002 | 14328.440 | UK pounds | UK pounds | UK | UK | Cheque | False | Normal_Small_Fan_Out |
| 4 | 10:35:21.0... | 2022-10-07 | 5376652437 | 9600420220 | 11895.000 | UK pounds | UK pounds | UK | UK | ACH | False | Normal_Fan_In |
| 5 | 10:35:21.0... | 2022-10-07 | 9614186178 | 3803336972 | 115.250 | UK pounds | UK pounds | UK | UK | Cash Deposit | False | Normal_Cash_Deposits |
| 6 | 10:35:21.0... | 2022-10-07 | 8974559268 | 3143547511 | 5130.990 | UK pounds | UK pounds | UK | UK | ACH | False | Normal_Group |
| 7 | 10:35:23.0... | 2022-10-07 | 980191499 | 8577635959 | 12176.520 | UK pounds | UK pounds | UK | UK | ACH | False | Normal_Small_Fan_Out |
| 8 | 10:35:23.0... | 2022-10-07 | 8057793308 | 9350896213 | 56.900 | UK pounds | UK pounds | UK | UK | Credit card | False | Normal_Small_Fan_Out |
| 9 | 10:35:26.0... | 2022-10-07 | 6116657264 | 656192169 | 4738.450 | UK pounds | UK pounds | UK | UK | Cheque | False | Normal_Fan_Out |
| 10 | 10:35:29.0... | 2022-10-07 | 7421451752 | 2755709071 | 5883.870 | Indian rupee | UK pounds | UK | UK | Credit card | False | Normal_Fan_Out |
| 11 | 10:35:31.0... | 2022-10-07 | 5119661534 | 9734073275 | 2342.310 | UK pounds | UK pounds | UK | UK | Debit card | False | Normal_Small_Fan_Out |
| 12 | 10:35:34.0... | 2022-10-07 | 5606024775 | 8646193759 | 1239.610 | UK pounds | UK pounds | UK | UK | Cash Deposit | False | Normal_Cash_Deposits |
| 13 | 10:35:34.0... | 2022-10-07 | 1405792899 | 5109623450 | 16555.310 | UK pounds | Pakistani rupee | UK | UK | Credit card | False | Normal_Fan_In |
| 14 | 10:35:37.0... | 2022-10-07 | 2188890133 | 3938416782 | 15459.460 | UK pounds | UK pounds | UK | UK | Cheque | False | Normal_Small_Fan_Out |
| 15 | 10:35:37.0... | 2022-10-07 | 6715177555 | 4460925916 | 586.280 | UK pounds | UK pounds | UK | UK | Cheque | False | Normal_Small_Fan_Out |

*Figure 2: SAML-D at a glance.*

An ongoing challenge in the AML framework is comparing the results of different Machine Learning algorithms, since the relative experiments are often conducted using datasets with distinct characteristics (Oztas et al., 2022). Thus, the main objective of the researchers who created the SAML-D dataset was to address this challenge by providing peer researchers with a challenging benchmark for evaluating classification models and enabling consistent results comparison, consequently supporting more meaningful analysis. Furthermore, SAML-D also aims to overcome the lack of data for AML analyses, mostly due to legal and privacy limitations that severely limit researchers' possibilities (Jullum et al., 2020).

## 4.5) Implementation and comparison of classification models

In this section, we describe the practical application of the models in the project, with particular attention to their general functioning and the reasons that guided their implementation. In this sense, after finding the best configuration, the main objective was to evaluate the performance of the Logic Learning Machine task in two different scenarios:

- **Pure Machine Learning:** In this setup, the Logic Learning Machine task is used considering exclusively the initial core attributes of SAML-D, without the integration of additional features based on heuristic rules (this type of models is called "Pure"). The

underlying idea was to set a sort of benchmark for this type of classification so that we could compare the results coming from other configurations of the LLM and check whether we obtained any improvements.

- **Combination of Machine Learning and heuristic:** In this case the Logic Learning Machine task is used on the "enriched" dataset we created by adding the features derived from the application of heuristic rules. The aim was to check and consequently gauge potential improvements in the accuracy and interpretability of the results, provided by the introduction of the heuristic rules. This procedure is carried out for two different model variants: the first containing all the heuristic rules we defined (the so-called "All" models), and the second one only containing the set of best rules we previously selected (the so-called "Best Rules" models). This was done to verify whether the Logic Learning Machine task worked better by using all the available information or only by providing it with a part of the total, i.e. the qualitatively better information.

In parallel to this analysis, we also proposed an additional study, using both the Decision tree and the Logistic classification tasks, following the same process proposed for the Logic Learning Machine. This was done primarily to obtain an accurate reference benchmark to compare the results and reach a more comprehensive view of the phenomenon.

In the integration phase between the Logic Learning Machine and the heuristic ruleset, the solution adopted consists in including the attributes derived from the heuristic rules among the input features, allowing the model to independently manage their information content. Subsequently, after the forecast phase, we manually intervene on the prediction, setting the predicted value as "fraudulent" every time a reasonable number of heuristic rules occurs, regardless of the result produced by LLM. This operation is handled by a module that assigns the value "1" every time one of the selected heuristic rules applies. If the sum of such applications exceeds a predefined threshold, the module sets the score value of the observation equal to "1"; otherwise, it retains its original value. Basically, we assign absolute priority to heuristic rules, considering their weight higher than those extracted by the model. However, this forcing process only occurs for a small subset of rules with an extremely low error ($\leq$0.001), which makes them almost flawless. This means that, although these rules rarely fire in the dataset, when they do, we consider them foolproof.

In parallel with this strategy, we also analyzed the binarization of the model. In fact, the standard cutoff of 0.5 for the score may not be optimal, especially in a fraud detection context with highly unbalanced classes. In these cases, the model generally tends to assign lower scores to most observations, leading the default threshold to ignore many frauds and generate a high number of false negatives. To improve the separation between classes, we therefore adopted a new threshold calculated in a data-driven way, exploiting the Youden index, which identifies the optimal balance point between recall and specificity. Once we obtained the new cutoff, we then re-binarized the model using this optimal threshold and computed the AUC metric and the confusion matrix to evaluate the models.

## 4.6) Self-coding development - Rulex Platform

The so-called self-code platforms are programs that combine the visual approach of no-code programs, i.e. a straightforward WYSIWYG (What You See Is What You Get) interface, with the possibility of developing complex projects thanks to the fact that the platform automatically writes the underlying code. The Rulex Platform is a clear example of a self-code platform, as it offers both a simple and intuitive drag and drop interface and optimized code for every operation it performs in its tasks.

Each executed action, such as data transformation, model creation and custom calculations, is traced and saved in an interactive history table and every step made can be saved, re-executed, undone or deleted. Each single operation can be inspected, making the code behind it visible. In the Rulex Platform, operations are carried out by linking together different blocks which perform specific operations following the direction of the flow. This allows to create a clear and organized succession of computations.

While being a low-code platform, the Rulex Platform offers many tools for implementing advanced customization, thanks to its built-in functions and formulas that can be directly configured by the user. Each individual function, formula and customizable parameter can be examined inside the specific task, allowing the user to access the relative documentation and better comprehend the tool being used. Another key feature of the Rulex Platform is that it includes Machine Learning tasks such as the Logic Learning Machine task, which allows to implement in depth analysis and evaluate models without the need of writing complex programming algorithms.

In short, the Rulex Platform is a versatile user-friendly tool that allows anyone to work with data and perform in-depth analysis, regardless of his initial set of programming skills. The combination of its intuitive interface and wide range of internal tasks makes it a valid tool both for business and academic contexts.

## 4.7) Flow in the Rulex Platform

In this section we briefly analyze the structure of the flow, describing its main components. A synthetic diagram of the flow is provided in Figure 3.
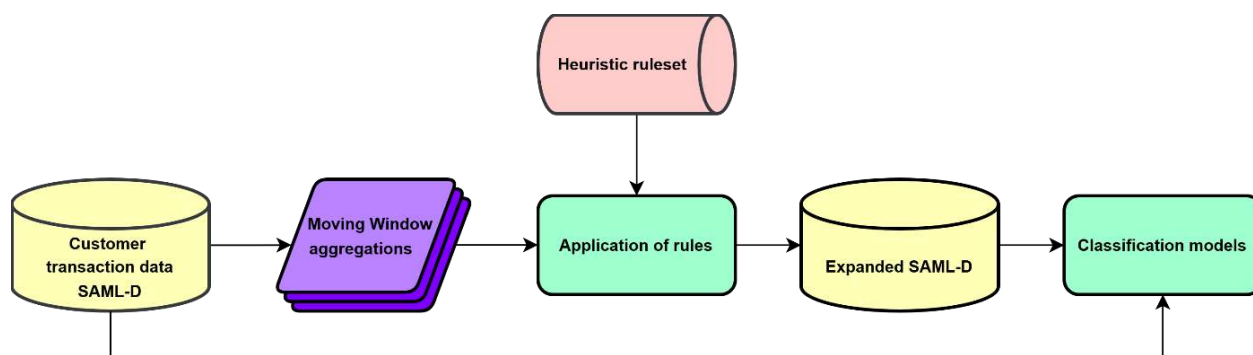


*Figure 3: Schematic flow diagram.*

After importing the dataset, we implemented the moving window aggregations. The moving window task performs data aggregations over defined time intervals, leveraging measures such as minimum, maximum, mean, and median. This operation can be applied to present, past, or future time frames, with the latter two achieved by shifting the window backward or forward by a fixed time delta. In this study, we employed both present and past aggregations to enable comparative analysis and identify potential suspicious differences in accounts' behaviors. The aggregated data generated through this process was stored in new attributes, which were subsequently used in the heuristic rule application phase. This step was conducted using the rule engine task in the Rulex Platform, which combines the information coming from a dataset with a predefined set of rules. The latter produces new attributes containing the flags raised by each individual rule, whenever a fraudulent pattern is detected.

The rule application phase served a dual purpose: first, to evaluate the heuristic ruleset by computing metrics such as covering, error, and their related performance indicators; second, to generate an expanded version of the original dataset, incorporating the outputs of the rule evaluation. This enhanced dataset was then used as input for the subsequent classification tasks. For consistency and to assess the added value of the heuristic ruleset, the same classification procedure was also applied to the original dataset, enabling a direct comparison of the results. A simplified representation of the flow developed in the Rulex Platform regarding the implementation of the three classification algorithms is provided in Figure 4. As described, the same was performed both on the original SAML dataset and on its enriched version containing heuristic.
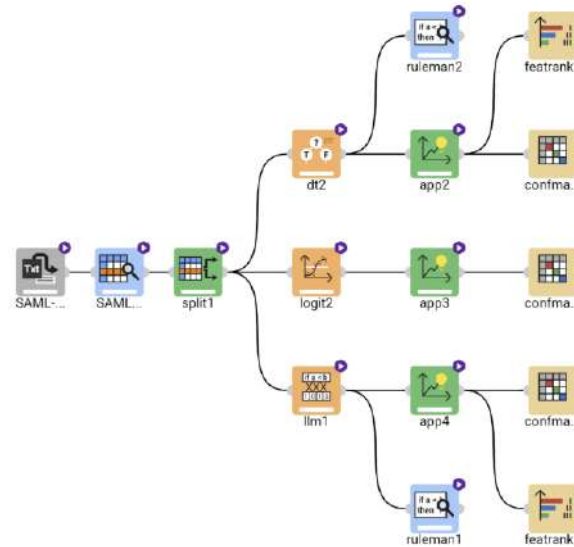


*Figure 4: Simplified representation of the Decision Tree, Logistic and Logic Learning Machine classification flow.*

## 5) Results

In this section, we present the results obtained in the various configurations of the models analyzed previously. Table 3 reports the performance of the baseline "Pure" models, providing a reference point for comparison with subsequent ones. In particular, the AUC values on the training and test sets (denoted with "$AUC_{training}$" and "$AUC_{test}$" respectively) and the computation times for the training and testing phases (denoted with "$t_{training}$" and "$t_{test}$" respectively and measured in seconds) are examined and compared.

|  | $AUC_{training}$ | $AUC_{test}$ | $t_{training}$ | $t_{test}$ |
|---|---|---|---|---|
| **LLM Pure** | 0.786 | 0.773 | 5125 | 205 |
| **Logistic Pure** | 0.761 | 0.770 | 297 | 325 |
| **DT Pure** | 0.687 | 0.677 | 627 | 15 |

*Table 3: "Pure" models results comparison.*

Table 4 shows the results of the models integrating the heuristic rules. The models marked with (*) indicate those in which the heuristic rules have been forced and the binarization has been updated, based on the optimal values obtained through the Youden index.

|  | $AUC_{training}$ | $AUC_{test}$ | $t_{training}$ | $t_{test}$ |
|---|---|---|---|---|
| **LLM All** | 0.861 | 0.841 | 19261 | 161 |
| **LLM All (*)** | 0.861 | 0.840 | 19261 | 161 |
| **Logistic All** | 0.830 | 0.828 | 524 | 306 |
| **DT All** | 0.687 | 0.677 | 635 | 16 |
| **LLM Best Rules** | 0.787 | 0.778 | 11281 | 124 |
| **LLM Best Rules (*)** | 0.787 | 0.778 | 11281 | 124 |
| **Logistic Best Rules** | 0.761 | 0.770 | 138 | 289 |
| **DT Best Rules** | 0.687 | 0.677 | 591 | 14 |

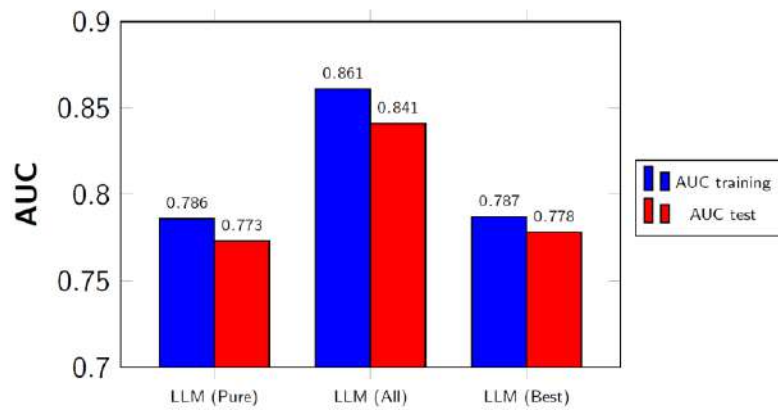*Table 4: "All", "Best rules" and "(*)" variants models results comparison.*

*Figure 5: LLM "Pure", "All Rules" and "Best Rules" models AUC comparison.*

In this context, it is also particularly useful to analyze ROC curves, as they represent the visual and conceptual counterpart of the AUC. Studying these two elements together provides a more complete view of the discriminative capacity of the model.

The AUC, in fact, makes it possible to assess the goodness-of-fit of the predictive score in a global manner, as it considers the performance of the model for each possible cutoff.

In contrast, point metrics such as the Youden J index focus on a specific cutoff, providing a 'snapshot' of performance at that point, but neglecting the overall behavior of the model.

The latter measures the vertical distance from the random classification line (the bisector starting at point (0,0)) and is therefore a useful tool to identify the optimal balance point between sensitivity and specificity.

In particular, maximizing it allows us to identify the most diagnostically effective decision threshold.

Figures 6, 7 and 8 show the ROC curves on the test set for the models presented in Figure 5.

The first thing that catches the eye when analyzing the results is the inferiority of the Decision tree models compared to the others, in terms of AUC.

However, this outcome was somewhat predictable, considering the intrinsic limitations of the methodology, particularly in contexts characterized by strong class imbalance, such as the one analyzed here.

In such settings, the Decision tree tends to perform poorly, as the algorithm is heavily influenced by the initial root-split, which in turn conditions all the subsequent splits, inevitably compromising the model's ability to effectively detect the minority class (i.e. fraudulent cases).

As a result, Decision tree models fail to produce truly informative trees, substantially returning the same output across all three configurations: "Pure", "All", and "Best Rules".

This highlights a clear limitation, making this methodology not particularly useful for a meaningful comparison with the other methods.

In short, given the significant performance gap, to allow a fairer and more reasonable comparison, the subsequent comparative analysis will focus primarily on the Logistic and Logic Learning Machine models.
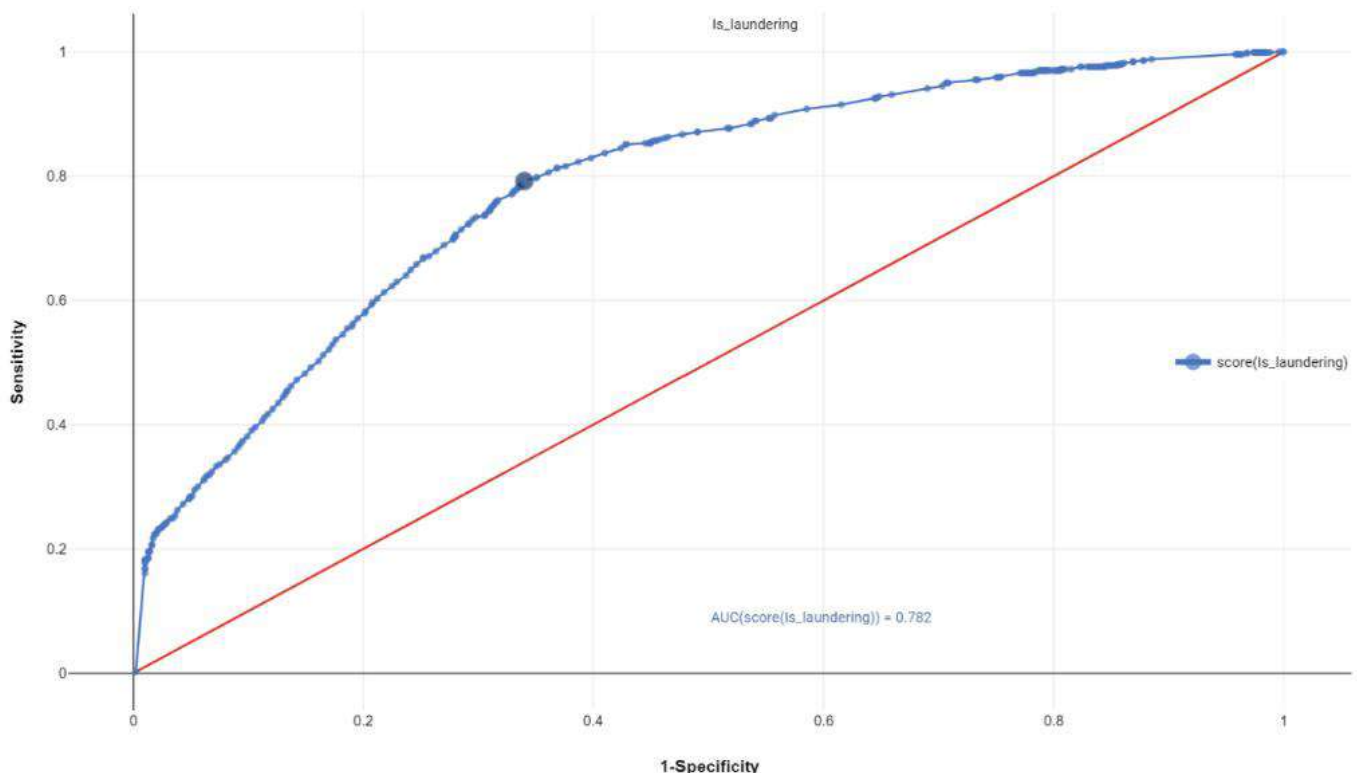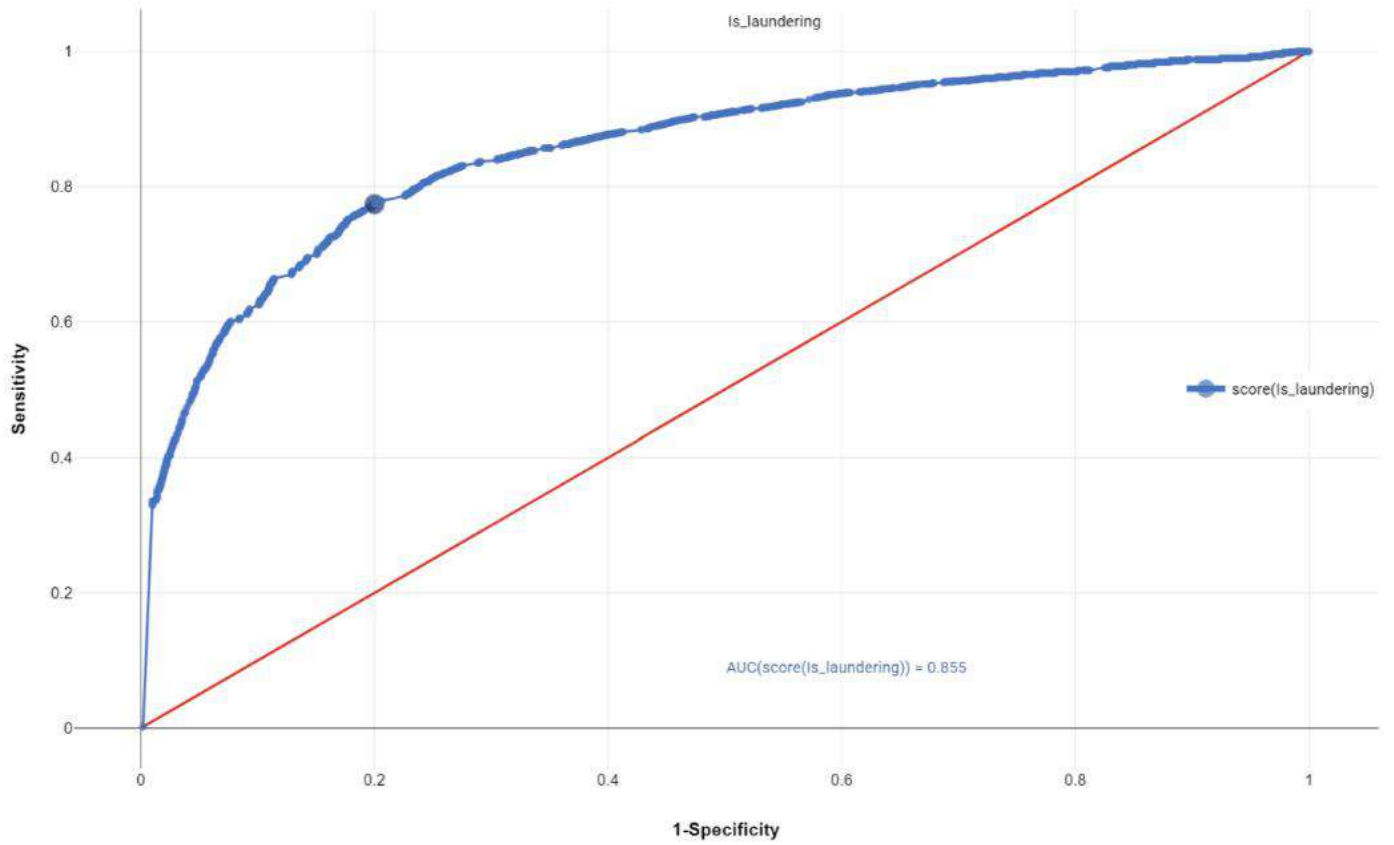


*Figure 6: LLM "Pure" ROC curve.*

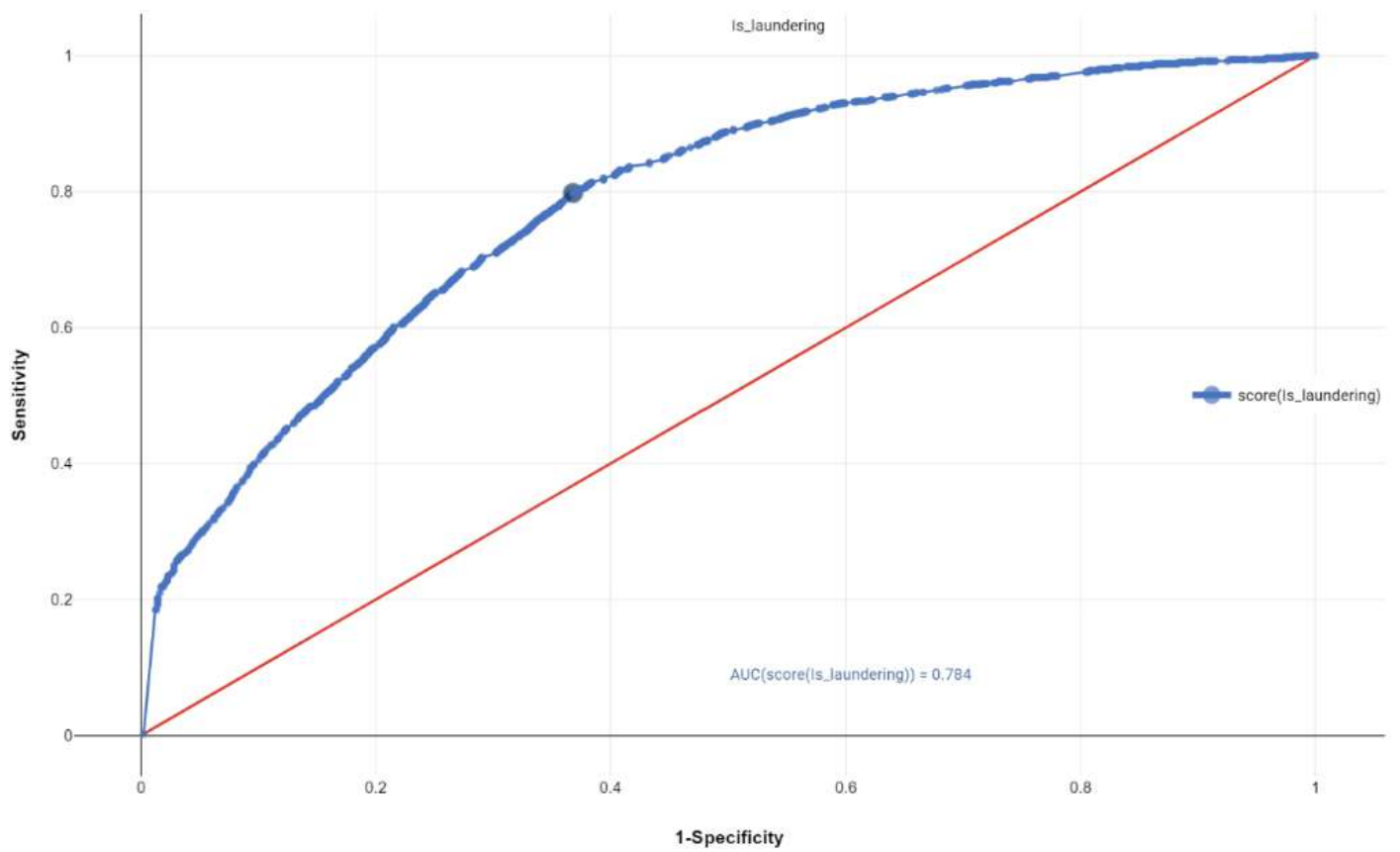*Figure 7: LLM "All" ROC curve.*



*Figure 8: LLM "Best Rules" ROC curve.*

## 5.1) Confusion matrices

To correctly interpret the results obtained from the confusion matrices, it is essential to remember the starting point of the original dataset. The latter, as already highlighted, presents a strong class imbalance, with only 0.1% of fraudulent cases.

Consequently, apparently low precision values are not necessarily disappointing, but instead they represent a significant improvement compared to random classification. For example, a precision of 1%, which may seem unsatisfactory at first sight, actually indicates a model ten times more precise than a random classification.

In general, the results obtained for precision and recall statistics tend to be extreme in opposite directions: either very low precision with high recall is recorded, or the opposite occurs. This imbalance makes both metrics uninformative for a significant descriptive analysis. For a more balanced evaluation, it is more appropriate to consider a synthetic indicator such as Youden's J statistic, which measures the discriminating capacity of the model, without privileging a single aspect, as precision and recall do. This allows for a more reliable evaluation of the overall effectiveness of the model, avoiding misleading interpretations due to unilaterally optimized metrics. A summary of the results is reported in Table 5:

|  | Youden's J statistic |
|---|---|
| **LLM Pure** | 0.151 |
| **Logistic Pure** | 0.005 |
| **LLM All** | 0.318 |
| **LLM All (*)** | 0.575 |
| **Logistic All** | 0.006 |
| **LLM Best Rules** | 0.173 |
| **LLM Best Rules (*)** | 0.430 |
| **Logistic Best Rules** | 0.005 |

*Table 5: "Pure", "All", "Best rules" and "(*)" variants Youden's J statistic results comparison.*

From the results reported in Table 5, it clearly emerges that the Logic Learning Machine models' configurations obtain the best performance. In particular, the LLM All (*) model stands out with a value of 0.575, showing a good discriminating capacity and demonstrating to be the most effective model in this setting, also considering the strong results previously observed in terms of AUC. In this case, the combination of the Logic Learning Machine and heuristic rules resulted in the generation of 243 rules. The best five rules, in terms of covering, are reported in Table 6. These rules are defined by the following conditions: as can be observed, many of these rules contain conditions regarding the transaction amount and payment types, suggesting the high informational value of these attributes for detecting suspicious behaviors. Some rules also contain conditions that exploit heuristic attributes regarding historical aggregation (e.g. "PayType1D1M" in rule #4), indicating that also the temporal evolution of transactions plays a relevant role in the identification of anomalies.

|  | Cond. #1 | Cond. #2 | Cond. #3 | Cond. #4 | Covering |
|---|---|---|---|---|---|
| **Rule #1** | "Amount" > 2033.765 | "Payment_currency" in [UK pounds] | "Payment_type" in [Cash Deposit] | "CurPairs1M1M H" in [No] | 11.206 |
| **Rule #2** | "Amount" > 2564.305 | "Payment_type" in [Cash Deposit] |  |  | 9.122 |
| **Rule #3** | "Payment_type" in [Cash Withdrawal] | "PayType1W1W H" in [No] |  |  | 8.664 |
| **Rule #4** | "Amount" > 177.865 | "Payment_type" in [Cash Deposit, Cash Withdrawal] | "PayType1D1M" in [Yes] |  | 7.732 |
| **Rule #5** | "Amount" <= 16434.655 | "Payment_currency" in [UK pounds] | "Payment_type" in [Cheque, Credit card, Debit card] | "PayType1M1M H" in [No] | 6.83 |

*Table 6: LLM "All (*)" first best rules in terms of covering.*

## 6) Conclusions

To choose the most suitable model, it is essential to adopt an objective-based approach, that is, to identify the most appropriate trade-off between two opposite scenarios, depending on the specific needs of the case. On the one hand, if the main objective is to identify fraudulent cases with maximum precision, it is appropriate to adopt models characterized by high levels of precision. However, this strategy involves an inevitable trade-off: a reduction in recall, with the risk of labeling many fraudulent cases as false negatives. This approach is particularly suitable when the system is unable to handle a large number of reports and must therefore prioritize the quality of identifications over quantity.

In real case scenarios, however, the main problem is not false negatives, which represent a necessary trade-off, but rather false positives, which can generate high costs and inefficiencies, without leading to an actual improvement in fraud detection.

On the other hand, if you want to preserve recall, i.e. maintain the maximum possible number of fraud reports, you need to opt for more balanced models. This choice allows you to intercept a greater number of illicit activities but also involves an increase in operational costs and potential inconvenience for customers. In our analysis, the model that has demonstrated the most balanced performance, and is therefore the most suitable choice for this scenario, is LLM All (*).

In short, the selection of the most appropriate model essentially depends on two factors that must be carefully balanced: the operational costs related to the management of reports and the number of frauds actually detected and reported.

This paper confirmed the results reported in the literature, demonstrating how the integration between Machine Learning and heuristic rules can significantly improve the predictive performance of classification models. In this specific case, with the adopted parameterizations, Logic Learning Machine models proved to be the best choice for fraud detection, clearly outperforming the Decision tree and Logistic regression algorithms. Among these, the most balanced and performing model was found to be LLM All (*), as it best combined the advantages of heuristic information with the benefits deriving from rule forcing and re-binarization.

Although the results obtained are very promising, there is still room for improvement. Specifically, a fundamental evolution concerns the continuous updating of the ruleset, introducing new rules to counter the evolution of money laundering techniques and guaranteeing a constantly effective detection system. Finally, a crucial step will be to replicate the analysis on new compatible transaction datasets. The money laundering problem is highly dependent on the quality and variety of available data, which implies that the performance of a model can vary significantly based on the information used for training and testing. Testing models on different datasets would allow to obtain more stable results and to conduct a more solid and coherent analysis.

Further developments could also explore the adoption of alternative model configurations specifically designed to maximize precision. while preserving an acceptable level of recall. By accurately analyzing and selecting the right value for the "confidence" parameter of the models, it is indeed possible to boost precision metric and control the trade-off with recall, tailoring the model's behavior to operational constraints. Moreover, the fine-tuning of heuristic rules through the Rule Enhancer module, a Rulex tool capable of automatically adjusting thresholds and nominal values based on a user-defined metric, may become a key tool for optimizing model performance according to specific goals. Finally, following recent calls for greater attention to safe machine learning (Giudici, 2024), future work could investigate the adoption of alternative model configurations specifically designed not only to enhance precision, but also to ensure robustness, transparency, and the integration of safety-oriented evaluation criteria.

## References

[1] Aparício, David and Barata, Ricardo and Bravo, João and Ascensão, João Tiago and Bizarro, Pedro (2020). Arms: Automated rules management system for fraud detection. In arXiv. Retrieved from: https://arxiv.org/abs/2002.06075 (accessed 13th June 2025).

[2] Berretta S., Fusaro M., Giribone P. G., Muselli M., Tropiano F., Verda D. (2025) – "Enhancing the explainability of the default probability model using the Logic Learning Machine: a comparison between native "white boxes" Machine Learning techniques" – International Journal of Financial Engineering. Online Ready. https://doi.org/10.1142/S2424786325500057.

[3] Chen, Zhiyuan and Van Khoa, Le Dinh and Teoh, Ee Na and Nazir, Amril and Karuppiah, Ettikan Kandasamy and Lam, Kim Sim (2018). Machine Learning techniques for anti-money laundering (AML) solutions in suspicious transaction detection: a review. Knowledge and Information Systems, 57, pp. 245-285.

[4] European Commission (2023). Anti-money laundering and countering the financing of terrorism at EU level. In finance.ec.europa.eu. Retrieved from: https://finance.ec.europa.eu/financial-crime/anti-money-laundering-and-countering-financing-terrorism-eu-level_en (accessed 13th June 2025).

[5] European Union (2018). Directive (EU) 2018/843 of the European Parliament and of the Council. In eur-lex.europa.eu. Retrieved from: https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:32018L0843 (accessed 13th June 2025).

[6a] European Union (2024). Regulation (EU) 2024/1620 of the European Parliament and of the Council. In eur-lex.europa.eu. Retrieved from: https://eur-lex.europa.eu/eli/reg/2024/1620/oj/eng (accessed 13th June 2025).

[6b] European Union (2024). Regulation (EU) 2024/1624 of the European Parliament and of the Council. In eur-ex.europa.eu. Retrieved from: https://eur-lex.europa.eu/eli/reg/2024/1624/oj/eng (accessed 13th June 2025).

[6c] European Union (2024). Regulation (EU) 2024/1640 of the European Parliament and of the Council. In eur-lex.europa.eu. Retrieved from: https://eur-lex.europa.eu/eli/dir/2024/1640/oj/eng (accessed 13th June 2025).

[7] Financial Action Task Force (FATF) (2025). The FATF Recommendations. In fatf-gafi.org. Retrieved from: https://www.fatf-gafi.org/content/dam/fatf-gafi/recommendations/FATF%20Recommendations%202012.pdf (accessed 20th November 2025).

[8] Financial Action Task Force (FATF) (2013). Targeted financial sanctions related to terrorism and terrorist financing (Recommendation 6). In fatf-gafi.org. Retrieved from: https://www.fatf-gafi.org/content/dam/fatf-gafi/guidance/BPP-Fin-Sanctions-TF-R6.pdf.coredownload.pdf (accessed 13th June 2025).

[9] Gaggero G., Giribone P. G., Muselli M., Ünal E., Verda D. (2024) – "Portfolio optimization and risk management through Hierarchical Risk Parity and Logic Learning Machine: a case study applied to the Turkish stock market – Risk Management Magazine Vol. 19, N. 1.

[10] Giudici, P. (2024). Safe machine learning. Statistics, Vol. 58(3), 473-477. https://doi.org/10.1080/02331888.2024.2361481

[11] International Revenue Service (IRS) (2025). Bank Secrecy Act. In irs.gov. Retrieved from: https://www.irs.gov/businesses/small-businesses-self-employed/bank-secrecy-act (accessed 13th June 2025).

[12] Jullum, Martin and Løland, Anders and Huseby, Ragnar Bang and Ånonsen, Geir and Lorentzen, Johannes (2020). Detecting money laundering transactions with machine learning. Journal of Money Laundering Control, Vol. 23 (1), pp. 173-186.

[13] KYC-CHAIN (2020). An overview of the FATF Recommendations, 2020. In kyc-chain.com. Retrieved from: https://kyc-chain.com/an-overview-of-the-fatf-recommendations/ (accessed 13th June 2025).

[14] Lewis, Roger J. (2000). An Introduction to Classification and Regression Tree (CART) Analysis. Conference proceedings of the Annual Meeting of the Society for Academic Medicine in San Francisco, California.

[15] London Stock Exchange Group (LSEG) (2024). EU Anti-money laundering directives. In lseg.com. Retrieved from: https://www.lseg.com/en/risk-intelligence/financial-crime-risk-management/eu-anti-money-laundering-directive (accessed 13th June 2025).

[16] Milo, Tova and Novgorodov, Slava and Tan, Wang-Chiew (2016). Rudolf: interactive rule refinement system for fraud detection. Proceedings of the VLDB, Vol. 9(13), pp. 1465-1468.

[17] Muselli, Marco (2005). Switching Neural Networks: A New Connectionist Model for Classification. Neural Nets, pp. 23-30.

[18] Muselli, Marco and Ferrari, Enrico (2011). Coupling Logical Analysis of Data and Shadow Clustering for Partially Defined Positive Boolean Function Reconstruction. IEEE Transactions on Knowledge and Data Engineering,Vol. 23, pp. 37-50.

[19] Nweze, Michael and Avickson, Eli and Ekechukwu, Gerald (2024). The Role of AI and Machine Learning in Fraud Detection: Enhancing Risk Management in Corporate Finance. International Journal of Research Publication and Reviews, Vol. 5.

[20] Ohno-Machado, Lucila and Stephan, Dreiseitl (2002). Logistic regression and artificial neural network classification models: a methodology review. Journal of Biomedical Informatics, Vol. 35(5–6), pp. 352-359.

[21] Oztas, Berkan and Cetinkaya, Deniz and Adedoyin, Festus and Budka, Marcin (2022). Enhancing Transaction Monitoring Controls to Detect Money Laundering Using Machine Learning. 2022 IEEE International Conference on e-Business Engineering (ICEBE), pp. 26-28.

[22] Oztas, Berkan and Cetinkaya, Deniz and Adedoyin, Festus and Budka, Marcin and Dogan, Huseyin and Aksu, Gokhan (2023). Enhancing Anti-Money Laundering: Development of a Synthetic Transaction Monitoring Dataset. 2023 IEEE International Conference on e-Business Engineering (ICEBE), pp. 47-54. Dataset Retrieved from: https://www.kaggle.com/datasets/berkanoztas/synthetic-transaction-monitoring-dataset-aml (accessed: 20 November 2025).

[23] Teradata. (n.d.). Fraud Detection with Machine Learning. In teradata.de. Retrieved from: https://www.teradata.de/insights/ai-and-machine-learning/fraud-detection-machine-learning (accessed 13th June 2025).

[24] UIF (Unità di Informazione Finanziaria per l'Italia) istruzioni per la rilevazione di operazioni sospette (2025). Retrieved from: https://uif.bancaditalia.it/normativa/norm-antiricic/Istruzioni_UIF_rilevazione_e_segnalazione_operazioni_sospette.pdf (accessed: 29th July 2025).

[25] United Nations (2021). Learn about UNCAC. In unodc.org. Retrieved from: https://www.unodc.org/corruption/en/uncac/learn-about-uncac.html (accessed 13th June 2025).

[26] United Nations Office on Drugs and Crime (UNODC) (n.d.). Overview of Money Laundering. In unodc.org. Retrieved from: https://www.unodc.org/unodc/en/money-laundering/overview.html (accessed 13th June 2025).

[27] United Nations Office on Drugs and Crime (UNODC) (2000). United Nations Convention Against Corruption. In unodc.org. Retrieved from: https://www.unodc.org/documents/treaties/UNCAC/Publications/Convention/08-50026_E.pdf (accessed 13th June 2025).

# Appendix A – Table of letters and symbols

| | |
|---|---|
| $Y \in \{O_0, O_1, \ldots, O_{q-1}\}$ | *Explanatory categorical variable with q classes (q= 2 is a bi-class problem)* |
| $q$ | *Number of classes of Y* |
| $O_0, O_1, \ldots, O_{q-1}$ | *Labels of the classes of Y* |
| $X = X_1, \ldots, X_j, \ldots, X_d$ | *Features* |
| $d$ | *Number of features* |
| $y_i , i = 1, \ldots, n$ | *Output of sample i* |
| $x_{ij}\ , i = 1, \ldots, N_{row}\ , j = 1, \ldots, d$ | *Dataset dimension [n, d]* |
| $S = \{(x_i, y_i)\}_{i=1}^{n}$ | *Training set* |
| $g(x)$ | *Model* |
| $\psi_j$ | *Mapping from continuous domain to discrete domain of j-th variable* |
| $M$ | *Number of discretization intervals* |
| $I_M = 1, \ldots, M$ | Set of positive integers up to M |
| $\boldsymbol{\gamma}_j = (\gamma_{j1}, \ldots, \gamma_{jm}, \ldots, \gamma_{jM_j-1})$ | $(M_j - 1)$ *cutoffs of variable* $X_j$ |
| $M_j$ | *Number of discretization intervals of variable $X_j$* |
| $\boldsymbol{\rho}_j = (p_{jl}, \quad \forall\, l = 1, \ldots, \alpha_j)$ | the vector of all the $\alpha_j$ values for input variable j in ascending order |
| $\alpha_j$ | *Number of distinct values of $X_j$* |
| $\varphi_j$ | Mapping for transformation of discretized domain into a binary domain |
| $u, v$ | $u < v$ if and only if $\varphi_j(u) < \varphi_j(v)$ |
| $\boldsymbol{z}, \boldsymbol{w} \in \{0,1\}^{M_j}$ | Elements of a string having a bit for each possible value in $I_{M_j}$ |
| $\boldsymbol{z}_i \in \{0,1\}^B$ ; $\varphi_j(\boldsymbol{x}_i) = \boldsymbol{z}_i$ | *Binarization of x* |
| $\boldsymbol{z}_i$ | Obtained by concatenating $\varphi_j(\mathbf{x}_i)$ for $j = 1, \ldots, d$. |
| $B = \sum_{j=1}^{d} M_j$ | Sum of the number of discretization intervals for $j = 1, \ldots, d$. |
| $S' = \{(\boldsymbol{z}_i, y_i)\}_{i=1}^{N}$ | *Binarized Training Set* |
| $N$ | *Number of rows in training set* |
| $\boldsymbol{h} \in \{0,1\}^{M_j}$ | Having all bits equal to 1 except the k-th bit which is set to 0 |
| $f$ | Boolean function |
| $T$ | T is the set containing $(\mathbf{z}_i, y_i)$ with $y_i = 1$ |
| $F$ | F is the set containing $(\mathbf{z}_i, y_i)$ with $y_i = 0$ |
| $\eta$ | Number of terms in definitions of AND OR |
| $\mathbf{a} \in \{0,1\}^B$ | Vector of zeros or ones of length "B" |
| $A \subset I_B$ | Antichain |
| $A^*$ | Simplified A , $A^* \subset A$ |
| $P(\boldsymbol{a})$ | The subset $I_B$ containing each i such that $a_i = 1$ |
| $U(\boldsymbol{a})$ | Upper shadow of $\mathbf{a}$, |
| $L(\boldsymbol{a})$ | Lower shadow of $\mathbf{a}$, |
| $\boldsymbol{h}_j$ | $\mathbf{z}$ was obtained by concatenating the results of the mapping $\varphi_j(\mathbf{x})$ and it can be split into substring $\mathbf{h}_j$ for each attribute |
| $V$ | The set of values associated with each $z_i$ |
| $r$ | *Rule* |
| $C(r)$ | *Covering of rule r* |
| $E(r)$ | *Error of rule r* |
| $TP(r), FP(r), TN(r), FN(r)$ | True positive, false positive, True negative, and false negative for rule $\mathbf{r}$ |
| $\varepsilon_{max}$ | Regularization parameter |
| $R(r)$ | Relevance |
| $\mathcal{R}$ | Complete ruleset |
| $\mathcal{R}_o = \{r \mid r \in \mathcal{R}, O(r) = o\}$ | Set of rules that predict class o |
| $\mathcal{R}_o^i = \{r \mid r \in \mathcal{R}, r \leq \boldsymbol{x}_i, O(r) = o\}$ | Set of rules satisfied by $\mathbf{x}_i$ that predict class o |

| | |
|---|---|
| $S(\boldsymbol{x}_i, o)$ | Score for each class o that measures how likely it is that $y_i = o$ |
| $P(o \mid \boldsymbol{x}_i)$ | Probability that $y_i = o$ given x_i |
| $\tilde{y}_i$ | The selected output is the one that maximizes the output probability |
| $c$ | Condition |
| $r'$ | Obtained by removing condition c from r |
| $J_{j1}, J_{j2}, \ldots, J_{M_j}$ | Intervals $J_{j1} = [-\infty, \gamma_{j1}]$, $J_{j2} = (\gamma_{j1}, \gamma_{j2}]$, etc. |
| $G_j = \{v_{j1}, v_{j2}, \ldots\}$ | The collection of the possible values assumed by $x_j$ |

# Appendix B – SAML-D dataset

## B.1) Core features

**Time**: the precise time of each individual transaction. It is standardized to "UCT+00:00" convention;

**Date**: it gives information about the transaction date. This, in addition to the **Time** feature, is an essential feature for tracking transaction chronology;

**Sender_account**: it contains the information about the sending account ID;

**Receiver_account**: it contains the information about the receiving account ID. In addition to the **Sender_account** feature, it helps uncover behavioural patterns and complex banking connections;

**Amount**: it indicates transaction values to identify suspicious activities. This value is already standardized to £ (UK pounds);

**Payment_currency**: it gives information about the currency used to make the payment;

**Received_currency**: it gives information about the currency received as payment. Both this and **Payment_currency** generally align with the location feature of the account, meaning that they conform with the prevalent currency of that specific country, but several mismatched instances are also present to add complexity;

**Sender_bank_location**: it contains information relating to the country from which money is sent;

**Receiver_bank_location**: it contains the information relating to the country where the money is being received. Together with **Sender_bank_location** it helps pinpointing high-risk regions for AML such as Mexico, Morocco and the UAE. The same account may carry out transactions from different countries, meaning that this information is not static across different instances;

**Payment_type**: it specifies the typology of settlement carried out by the sender account, each involving different levels of risk. It includes various methods like credit card, debit card, cash, ACH transfers, cross-border, and cheque;

**Is_laundering**: this feature is a binary indicator differentiating "normal" from "suspicious" transactions;

**Laundering_type**: it further describes the typology of the transaction, classifying both "normal" and "suspicious" instances. It offers deeper insights into prevalent or high-risk typologies of transactions.

## B.2) Payment typologies

**Credit card;**

**Debit card;**

**Cash withdrawal;**

**Cash deposit;**

**Automated clearing house (ACH) transfers;**

**Cross-border;**

**Cheque.**

## B.3) Laundering typologies

**Cash withdrawal:** it involves withdrawing illicit funds in cash from a financial institution. It is used to move money out of the formal financial system and into the physical world, making it more difficult to trace. Cash withdrawals are often part of larger schemes, such as layering or integration in the money laundering process;

**Behavioural change 1:** the behavioural change 1 and 2 typologies adopt the same structure as the normal group typology. However, in behavioural change 1, the main account deviates from its usual patterns and transacts with new accounts;

**Behavioural change 2:** in contrast, under the Behavioural Change 2 typology, the main account transacts with new accounts in high-risk locations;

**Structuring:** it involves breaking large amounts of money into smaller transactions that fall below reporting thresholds;

**Smurfing**: this is a particular form of structuring where multiple individuals (called "smurfs") are employed to make numerous small deposits or transactions that are below the reporting thresholds;

**Fan out: w**hen a single source of illicit funds is distributed to multiple accounts or entities;

**Fan in:** when multiple smaller amounts from different accounts or entities are funneled into a single account;

**Layered fan in:** it is a multi-layer scheme which involves multiple sender accounts transacting with a fewer number of receiver accounts, which ultimately transact with a single receiver account in a sort of funnel-shaped pattern;

**Layered fan out:** this version mirrors the previous structure but with transactions flowing in the reverse direction;

**Scatter-gather:** when funds are first scattered across multiple accounts or entities (scatter phase) and then later reassembled into one or more accounts (gather phase);

**Gather-scatter:** it works in the exact opposite way;

**Cycle:** when funds are moved in a circular manner through a series of transactions that ultimately return the funds to the original account or entity;

**Bipartite:** when funds are transferred back and forth between two distinct groups of accounts in a manner that disguises the money's origin and destination;

**Stacked bipartite:** in its "stacked" version, multiple layers or tiers of entities are involved;

**Over-invoicing:** this technique exploits commercial transactions, e.g. transfer prices, to overestimate the value of goods and services provided to foreign affiliated partners in order to shift the taxable income to high-tax or low-tax jurisdictions;

**Deposit-send:** this typology refers to a situation where an account first deposits cash into the bank and then within a short period of time sends it to another account. The transaction amount is generally below the reporting threshold limit, with the second transaction having an increased chance of being sent to a high-risk country. It is considered suspicious due to the rapid movement of funds and potentially facilitating terrorism finance.

**Single large transaction:** as the name suggests, this type of payment involves sending a large sum of money in a single installment. This type of payment appears to be suspicious, especially in cases where a customer typically makes modest transactions and suddenly completes a single large transaction, which can be a sign of suspicious activity, particularly if there is no plausible economic justification.

# Examining cointegration between corporate governance and financial performance in selected listed South African financial institutions

Floyd Khoza (University of Mpumalanga - South Africa), Patricia Lindelwa Makoni (University of South Africa – UNISA)

Corresponding author: Patricia Lindelwa Makoni (makonpl@unisa.ac.za)

## Abstract

Since the 2007 global financial crisis, scholars have attempted to explain market failures using aspects other than corporate governance. Previous studies focused on the function of corporate governance in financial performance, with a dearth of literature on other financial dimensions like risk appetite, financial stability, and the effect of financial performance on the corporate governance of financial institutions. This study examines the cointegrating relationships between financial performance and corporate governance in selected South African listed financial institutions between 2007 and 2020. Employing the pooled mean group and fixed dynamic effect estimators in a panel autoregressive distributed lags framework, our results indicated notably positive long-term cointegrating relationships between the capital adequacy ratio (CAR) and financial stability, return on equity (ROE) and return on assets (ROA), respectively. We consider this paper valuable in that it contributes to the literature on the interrelatedness of corporate governance and financial performance, particularly of listed financial institutions, and is useful to central banks, market regulators, boards of directors and academics to inform policies and regulations.

## 1. Introduction

Several cases of institutional failures or collapses have been witnessed in the financial sector. United States (US) Financial institutions like Lehman Brothers, Washington Mutual, Wachovia, IndyMac Bank and J.P. Morgan (Nyaloti, 2024). International disasters in financial failures included non-financial firms such as the Maxwell saga in the United Kingdom (UK), Parmalat in Italy, Daewoo in Korea, and Macmed and Sentula in South Africa. Nigerian financial institutions included the Oceanic Bank, Savannah Bank Plc and Bank of the North (Gwaison and Maimako, 2021). In South Africa, financial sector scandals included Regal Treasury Bank, African Bank Saambou, Leisurenet, Fidentia, Venda Building Society (VBS) Mutual Bank, and JCI, demonstrating the growing need for transparency and robustness in governing the financial firms. Furthermore, South Africa reported on management misconduct in advisory firms such as Deloitte, African Bank and Klynveld Peat Marwick Goerdeler (KPMG) scandals (Lingwati and Mamabolo, 2023).

With the collapse of financial institutions and the activities of some other institutions, concerns have been raised about the need to improve corporate governance in financial institutions. According to Hunjra *et al*. (2024), sound corporate governance will ensure the effective and efficient functioning of financial institutions and the banking sector. Karpoff (2021) considers corporate governance to be an array of control procedures that organisations implement to restrict or discourage potentially self-interested managers (agents) from engaging in behaviours that are unfavourable to the financial welfare of shareholders and other stakeholders. Corporate governance describes how managers in charge of the company should run it. Therefore, the importance of the board of directors in institutionalising effective corporate governance principles in every organisation cannot be overstated.

The importance of the board of directors in corporate governance is apparent in model definitions of corporate governance, which define corporate governance as the structures and processes through which an institution's operations are directed and managed to improve long-term shareholders' value by improving corporate performance and accountability while taking stakeholders' interests into account (Tricker and Tricker, 2015).

The 2007 to 2009 global financial crisis emerged from corporate governance failures in the financial sector. Against this background, this study assesses the cointegrating relationships between financial performance and corporate governance in selected financial institutions.

## 2 Literature review
### 2.1 Theoretical literature review

The study was centred on both agency and stewardship theories' perceptions to enhance the impact of corporate governance on companies. According to Brealey *et al*. (2022) and Efunniyi *et al*. (2024), accounting information and corporate governance procedures can present stakeholders with information about an institution's financial position and performance. Accounting information summarises financial data in the form of ratios as the basis for forecasting future financial performance, which shareholders can use to make investment decisions (Brealey *et al*., 2022). Following the agency relationship, businesses experience agency problems due to asymmetric information between management, who act as agents, and shareholders, who act as principals involved in decision-making. As a result of random disturbances in the outcome of their actions, such as inefficient behaviour of all parties (i.e., shareholders and managers) in satisfying their interests, this information asymmetry may result in an incomplete contract (Schroeck, 2002). In agency theory, an agent is vital in formulating firm policies to provide investors with a sign and a desirable investment signal. However, good corporate development will inform stakeholders that the business has been able to maintain and grow its viability.

According to the stewardship theory, directors can accomplish organisational objectives for shareholders by improving their worth rather than self-serving. Donaldson and Davis (1991) support the argument of stewardship theory. According to stewardship theory, allowing managers (agents) to act with discretion can motivate them to perform better. According to Donaldson and Davis (1991), the stewardship approach emphasises that managers' concern for their career progression and reputation motivates them to operate in the best interests of the shareholders, reducing agency costs.

## 2.2 Empirical literature review

Oladipupo and Kelvin (2024) examined the impact of corporate governance on the financial performance of 39 listed Nigerian manufacturing firms from 2003 to 2022, in which the panel regression technique was employed. The results found that board size, audit committee and board composition negatively correlated with ROE. However, the audit committee and board composition positively correlated with ROA. The study found a negative relationship between board size and ROA. Furthermore, a positive relationship between board independence and financial performance (ROE and ROA). An earlier study by Elbahar (2016) assessed the key concepts of corporate governance, bank risk and performance of 90 banks selected from the GCC. Applying variables such as the return on assets (ROA) and the return on equity (ROE), he concluded that corporate governance variables of board size, gender diversity, duality and audit committee exerted little influence on bank performance (ROE). The study employed the Ordinary Least Squares (OLS) to analyse the results. Owusu and Garr (2024) investigated the impact of corporate governance on the financial performance of 14 Ghanaian-listed banks from 2008 to 2020. The study employed the Ordinary Least Squares (OLS) regression and found that board diversity, audit committee, board independence and board size had a positive relationship with ROA. However, a significant negative association was found between gender diversity and ROE.

Simanjuntak and Alfredo (2024) investigated the impact of corporate governance on the financial performance of Indonesian companies. The study employed regression analysis and sampled 100 publicly listed financial and non-financial companies. Shareholder rights, audit committee effectiveness and board independence proxied corporate governance, while ROA, net profit margin and ROE measured performance. The study found a positive relationship between corporate governance and financial performance. Talatu (2024) examined the effects of corporate governance on the financial performance of quoted healthcare firms in Nigeria from 2014 to 2023. The study employed the OLS regression model and found that board independence, skilled and diverse board composition positively correlate with ROE. Consistent with Mahmudi (2024), who found a positive correlation between financial performance and corporate governance using a systematic literature review.

Musa (2020) examined the association between corporate governance and the financial performance of Nigerian banks from the period 2013 to 2015. The study's independent variables are board meetings, board independence, and board gender, whereas the dependent variable is the ROA. The study sampled 15 listed banks and employed panel regression analysis. The results found the link between board independence and ROA is statistically insignificant. ROA and board meetings were found to be negatively significant. The association between board genders, the board size, and ROA, on the other hand, was statistically insignificant. Meanwhile, a positive and statistically significant association exists between firm size and ROA. The association between bank age and ROA was statistically significant and negative.

Kiptoo *et al*. (2021) examined the relationship between corporate governance and the financial performance of insurance firms in Kenya using data drawn from 51 Kenyan insurance companies from 2013–2018. Their corporate governance was measured exclusively using the board of directors' structure within the respective firms. They concluded that smaller boards, as well as boards with greater independent directors, were more efficient in enhancing financial performance within the surveyed insurance companies. Their findings were corroborated by Temba *et al*. (2023) in their assessment of the moderating role of corporate governance on the financial performance of commercial banks in Tanzania. Similarly, they recommended that enhanced corporate governance standards could improve the commercial banks' financial performance, particularly those pertaining to liquidity, capital adequacy, earning ability, efficient use of equity and asset quality.

Jouirou and Jouini (2022) examined the effect of gender diversity and directors' independence on French banks' performance. The study employed panel data regression model for 66 sampled French banks from 2015 to 2019. The results found a significant and positive relationship between gender diversity and profitability. In addition, a significant and positive effect of the independence of directors on the profitability of banks was found. However, Nizam and Liaqat (2022) examined the effect of corporate governance factors on bank performance in Pakistan, employing a sample of 15 banks from 2010 to 2020, using data from the financial reports. Their study applied the cointegration test, the Hausman test to determine the fixed or random effects and the Panel least squares regression to check the association between the variables. Nizam and Liaqat (2022) found a positive and significant association between board size and ROA, stipulating that optimal board size improves the ROA. In addition, their findings showed that board independence significantly affects ROA, implying that it plays a role in increasing shareholders' value.

Usendok (2022) investigated the relationship between corporate governance and firm performance in the Nigerian banking sector. ROA measures the firm performance of banks. Ex-post facto research and descriptive design were adopted, and the indirect least squares were employed in the study. The study sampled banks from 2014 to 2020, and the results found a significant positive association between board composition and performance, while board size and board meetings had a significant and negative relationship with the firm performance of the banks. Msomi and Nzama (2023) sought to identify firm-specific variables that affect the financial performance of 36 publicly listed South African insurance companies. Applying ROA as the dependent variable, they found that only leverage and liquidity ratios positively influenced financial performance, implying that insurance companies should focus on improving and maintaining these aspects. Muzata and Marozva (2023) assessed the effect of corporate governance on the financial performance of the Top 40 listed companies on the South African Johannesburg Stock Exchange (JSE). Based on their findings, they recommended that companies prioritise sound corporate governance practices because they positively influence firm performance.

**Table 1: Summary of literature review**

| Author's name and year | Title | Methodology | Findings |
|---|---|---|---|
| Oladipupo and Kelvin (2024) | Corporate governance and manufacturing firms' financial performance in Nigeria | Panel regression technique | The study found that board size, audit committee and board composition have a negative relationship with ROE. Audit committee and board composition had a positive relationship with ROA. A positive relationship between board independence and financial performance. |
| Elbahar (2016) | Corporate governance, risk management and bank performance in the GCC banking sector | OLS regression analysis | The study found board size, gender diversity, duality and audit committee exerted little influence on ROE. |
| Owusu and Garr (2024) | Corporate governance dynamics and financial performance: Analysis of listed commercial banks in the Ghanaian context | OLS regression analysis | The study found that board diversity, audit committee, board independence and board size had a positive relationship with ROA. |
| Simanjuntak and Alfredo (2024) | impact of corporate governance on financial performance Evidence from Indonesia. | Regression analysis | The study found a positive relationship between corporate governance and financial performance. |
| Talatu (2024) | Effect of corporate governance on financial performance of quoted healthcare firms in Nigeria. | OLS regression model | The study found that board independence, skilled and diverse board composition, positively correlated with ROE. |
| Mahmudi (2024) | Corporate governance mechanisms and financial performance: A systematic literature review in emerging markets | Systematic literature review | The study found a positive correlation between financial performance and corporate governance |
| Musa (2020) | Corporate governance and financial performance of Nigeria listed banks. | Panel regression analysis | The study found the link between board independence and ROA is statistically insignificant. ROA and board meetings were found to be negatively significant. The association between board genders, the board size, and ROA, on the other hand, was statistically insignificant. Meanwhile, a positive and statistically significant association exists between firm size and ROA |

| | | | |
|---|---|---|---|
| Kiptoo et al. (2021) | Corporate governance and financial performance of insurance firms in Kenya. | Regression analysis | The study found that smaller boards, as well as boards with greater independent directors, were more efficient in enhancing financial performance. |
| Jouirou and Jouini (2022) | Corporate governance mechanisms and banking performance. | Panel data regression model | The study found a significant and positive relationship between gender diversity and profitability. In addition, a significant and positive effect of the independence of directors on the profitability |
| Nizam and Liaqat (2022) | Corporate governance and firm performance: Empirical evidence from Pakistan banking sector | Panel least squares regression analysis | The study found a positive and significant association between board size and ROA, and showed that board independence significantly affects ROA. |
| Usendok (2022) | Corporate Governance and Organizational Performance: A Study of Selected Banks in Nigeria. | Indirect least squares | The study found a significant positive association between board composition and performance, while board size and board meetings had a significant and negative relationship with the firm performance |
| Muzata and Marozva (2023) | The Nexus between Executive Compensation and Firm Performance: Does Governance and Inequality Matter? | Generalised method of moments model analysis | They recommended that companies prioritise sound corporate governance practices because they positively influence firm performance. |

Source: Authors' own composition

Although many studies on corporate governance have been carried out across parts of the world, most of these studies have been conducted in different countries, employing different methodologies. South Africa has a robust regulatory framework for financial institutions. Furthermore, based on the scope of the research, limited studies have evaluated the impact of corporate governance in South Africa, utilising current data and different combinations of corporate governance measures employed in this study. This investigation was necessitated by the critical gap in the literature.

## 3 Data and methodology

Similar to the work of Khoza *et al.* (2024), the Bureau Van Dijk Orbis Bank and the Financial Sector Conduct Authority (FSCA) databases were used to source data on the financial institutions for this study. Although the South African financial services sector is fairly developed, our sample was restricted to 11 commercial banks and 10 insurance companies with complete data for the period under review (2007-2020).

Diagnostic tests were performed before data analysis to prevent spurious regression results. Principal component analysis (PCA) was used in this study to create a composite indicator of corporate governance. This method was adopted because there is no agreement in the literature on the most relevant variable to measure corporate governance (Swedan and Ahmed, 2019). To perform PCA, the Eigenvalues of the variance matrix must be computed. Several mutually independent principal components are applied to summarise the variables of interest, with each principal component being the weighted average of the underlying variables (Greenacre *et al.*, 2022). Meanwhile, the composite index constructed provides a methodologically efficient approach to lower dimensionality while capturing the overall governance quality. However, it may introduce limitations that can influence certain results. The PCA assumptions of linearity and orthogonality could not adequately capture the complexity and interactive nature of corporate governance measures. Therefore, it can conceal the opposing effects, which may result in an overall insignificant association with financial performance.

Board size was measured by the total number on the board of directors, independent non-executive directors measured by the total number of independent non-executive directors to total non-executive directors, non-executive directors measured by the number of non-executive directors to the total number of directors, board remuneration measured by the total amount of remunerations for

board members, board diversity measured by the percentage of female board members on total board members. Transparency and disclosure are measured by disclosures of financial information, general corporate governance disclosure, board of directors' reports, age and qualification of directors, compliance reports, committees, accounting policies, remuneration of directors, and auditors' reports. Therefore, the study employed PCA to combine the six corporate governance metrics into a single index, GOVINDEX. This study employed the corporate governance index as the dependent variable. The independent variables were risk appetite, financial stability, and financial performance. Financial performance is measured by return on assets (net income to average total assets) and return on equity (net income to average total equity). The Z-score measures financial stability (Kajumbula and Makoni, 2024), while risk appetite is proxied by CAR and measured by the capitalisation ratio, consisting of total equity to total assets. Biresaw and Sibindi (2025) confirmed that risk appetite was a necessary variable to measure and account for in financial institutions, as it contributes to their overall enterprise risk management (ERM) framework.

## 3.1 Panel autoregressive distributed lags

Pesaran *et al*. (1999) introduced the pooled mean group (PMG), dynamic fixed effects (DFE), and the mean group (MG) approaches in the Autoregressive Distributive Lags (ARDL) framework, which the current study followed. The PMG has the advantage of allowing financial institutions' heterogeneity in error variances, short-run coefficients, and intercepts, as well as the speeds of adjustments to long-run equilibrium values with a proposal of homogenous long-run slope coefficients across financial institutions (Pesaran *et al*., 1999). The MG estimator requires a separate equation for each cross-sectional dimension, and the model's parameters are averaged to create reliable estimators. The DFE presupposes that long-run coefficients are constant throughout the sample. The Hausman test determined the most suitable estimation technique among the MG, PMG and DFE. A dynamic model is preferred because corporate governance is persistent. This study jointly estimates the long-run and short-run impacts by employing the ARDL and the error correction model (ECM) in panel data.

Baltagi (2008) and Croissant and Millo (2019) argued that panel data presume heterogeneity, in contrast with either cross-sectional (N) or time series (T) studies. When heterogeneity is ignored, when individual institution-specific variables are not controlled, a model is mis-specified (Baltagi, 2008). Panel data improves the effectiveness of econometric estimates by providing the researcher with an increasing degree of freedom, a larger number of data points and decreasing multicollinearity across study variables (Hsiao *et al*., 1995; Fujiki *et al*., 2002; Baltagi, 2008; Hsiao, 2022). Furthermore, panel data allow a researcher to use aggregate data to examine critical economic problems not addressed with time series or cross-sectional data sets (Baltagi, 2008; Hsiao, 2022).

The financial dimension is the determinant of corporate governance in this study. Corporate governance is hypothesised to be the function of the financial dimension, namely, financial performance, risk appetite and financial stability.

The unrestricted panel ARDL is specified below:

$$\text{GOV}_{it} = \varphi_0 + \sum_{j=1}^{p} \delta_{it} \text{GOV}_{i,\,t-j} + \sum_{j=0}^{q} \delta_{2t} X_{i,t-j} + \mu_i + \varepsilon_{it} \quad \text{Eq. 1}$$

Where:

GOV is the dependent variable captured in this study, which includes board size (BS), board remuneration (BR), board diversity (BD), and board composition (BC). The corporate governance (GOV) proxies are regressed individually. k is the selected financial institution, with lag lengths p and q, respectively. $X_{i,\,t-j}$ is the vector of the explanatory variables for group i. $\mu_i$ is the fixed effect. $\varepsilon_{it}$ is the error term.

The equations below are re-parameterised to the specifics of the current study.

$$\text{GOV}_{it} = \beta_0 + \beta_{1i}\text{GOV}_{i,t-1} + \beta_{2i}\text{FINPERF}_{i,t-1} + \beta_{4i}\text{FINSTAB}_{i,t-1} + \sum_{j=1}^{p} \delta\Delta\text{GOV}_{i,t-j} + \sum_{j=0}^{q} \delta_{2t}\Delta\text{FINPERF}_{i,t-j} + \sum_{j=0}^{q} \delta_{4t}\Delta\text{FINSTAB}_{i,t-j} + \varepsilon_{it} \quad \text{Eq. 2}$$

$$\text{FINPERF}_{it} = \beta_0 + \beta_{1i}\text{FINPERF}_{i,t-1} + \beta_{2i}\text{GOV}_{i,t-1} + \beta_{4i}\text{FINSTAB}_{i,t-1} + \sum_{j=1}^{p} \lambda_{1t}\Delta\text{FINPERF}_{i,t-j} + \sum_{j=0}^{q} \delta_{2t}\Delta\text{GOV}_{i,t-j} + \sum_{j=0}^{q} \delta_{4t}\Delta\text{FINSTAB}_{i,t-j} + \varepsilon_{it} \quad \text{Eq. 3}$$

$$\text{FINSTAB}_{it} = \beta_0 + \beta_{1i}\text{FINSTAB}_{i,t-1} + \beta_{2i}\text{GOV}_{i,t-1} + \beta_{4i}\text{FINPERF}_{i,t-1} + \sum_{j=1}^{p} \lambda_{1t}\Delta\text{FINSTAB}_{i,t-j} + \sum_{j=0}^{q} \lambda_{2t}\Delta\text{GOV}_{i,t-j} + \sum_{j=0}^{q} \lambda_{4t}\Delta\text{FINPERF}_{i,t-j} + \varepsilon_{it} \quad \text{Eq. 4}$$

Where:

GOV is the corporate governance proxy, namely, board size (BS), board remuneration (BR), board diversity (BD) and board composition (BC) regressed individually. The proxies are regressed individually for corporate governance (GOV). FINPERF is the financial performance proxied by ROA and ROE. β denotes the long-run coefficient of the independent variable. Financial stability is FINSTAB. The short-run coefficients are φ, δ, γ, λ, Θ. The Akaike information criterion is applied to select the lag order (p, q). *t-j* represents the short-run and long-run relationships tested with differenced and lagged variables of the ARDL. The error term with the *i* of the institution and time period *t* is $\varepsilon_{it}$.

## 3.2 Error correction model

Once the long-run relationship between corporate governance and financial performance has been established, this study employs the vector error correction model (VECM) to determine the short-run effects (Apergis and Payne, 2010; Animasaun *et al.*, 2025). Engle and Yoo (1987), Phillips (1991), and Hoffman and Rasche (1996) argue that the error correction model provides the advantages of incorporating cointegrations and capturing the short-run effect of the variables being analysed. However, VECM uses the error correction model (ECM). ECM is employed instead of VECM if there is no cointegration.

The generic error correction model is therefore specified below:

$$\Delta GOV_{i,t} = \alpha_{0,t} + \sum_{j=1}^{p} \beta_j \Delta GOV_{i,t-j} + \sum_{j=0}^{q} \phi_{i,j} \Delta X_{i,t-j} + \varphi_{1i} ECT_{i,\,t-1} + \omega_{it} \quad \text{Eq. 5}$$

Where:

$\Delta$ denotes the first difference operator. GOV represents each of the corporate governance proxies, board size (BS), board remuneration (BR), board diversity (BD), and board composition (BC), which are regressed individually. B and $\phi$ denote the short-run coefficients. ECT represents the error correction term. X denotes the vector of the independent variable. $\varphi$ denotes the speeds of adjustments to the long-run equilibrium. $\alpha$ represents the constant. (p, q) represents the lagged lengths selected using the AIC. $\omega$ denotes the error term, which assumes a normal distribution with constant variance and zero mean.

The ECT coefficient ($\varphi$) in the ECM equation specifies the speed of system adjustments to long-run equilibrium after shocks in the short run. The ECT coefficients are expected to be statistically significant and negative to show convergence of the variables to equilibrium (Croissant and Millo, 2019).

We used the GOVINDEX as a measure of the corporate governance index. Equations 6 to 8 are equations for vector error correlation among corporate governance proxied by GOVINDEX (BS, BC, BC, BR, BD) and the financial variables (financial performance and financial stability). The equations are specified as follows:

$$\Delta GOVINDEX_{it} = \alpha_0 + \sum_{j=1}^{p} \delta \Delta GOVINDEX_{i,t-j} + \sum_{j=0}^{q} \delta_{2t} \Delta FINPERF_{i,t-j} + \sum_{j=0}^{q} \delta_{4t} \Delta FINSTAB_{i,t-j} + \phi_{1i} ECT_{i,\,t-1} + \varepsilon_{it} \quad \text{Eq. 6}$$

$$\Delta FINPERF_{it} = \alpha_0 + \sum_{j=1}^{p} \delta \Delta FINPERF_{i,t-j} + \sum_{j=0}^{q} \delta_{2t} \Delta GOVINDEX_{i,t-j} + \sum_{j=0}^{q} \delta_{4t} \Delta FINSTAB_{i,t-j} + \phi_{1i} ECT_{i,\,t-1} + \varepsilon_{it} \quad \text{Eq. 7}$$

$$\Delta FINSTAB_{it} = \alpha_0 + \sum_{j=1}^{p} \delta \Delta FINSTAB_{i,t-j} + \sum_{j=0}^{q} \delta_{2t} \Delta FINPERF_{i,t-j} + \sum_{j=0}^{q} \delta_{4t} \Delta GOVINDEX_{i,t-j} + \phi_{1i} ECT_{i,\,t-1} + \varepsilon_{it} \quad \text{Eq. 8}$$

Where:

GOVINDEX represents the corporate governance proxies: board size (BS), board remuneration (BR), board diversity (BD) and board composition (BC). GOVINDEX is a PCA composite index from the four individual proxies. FINPERF denotes the financial performance proxies, which are ROA and ROE. FINSTAB represents financial stability. $\varphi$, $\phi$, $\lambda$ denote speeds of adjustment to the long-run equilibrium. $\alpha$ denotes the constant. $\beta$ denotes the coefficients in the short run.

## 4 Results and discussion of the findings

This section presents the results and discussions of the cointegrations and the ECT between the corporate governance index and financial variables: financial performance, risk appetite and financial stability. Each financial variable is jointly assessed to investigate the cointegration relationship with the corporate governance index. Using the Hausman tests, the PMG and DFE are preferred estimators for the study as the coefficients are verified for long-run homogeneity. Section 4.1 discusses the panel cointegrating and the error correction model results based on the financial performance measure (ROA). Section 4.2 discusses the panel cointegrating and the error correction model results based on the financial performance measure (ROE). Tables 1 to 8 provide the results of the preferred estimators.

## 4.1 Panel Cointegration and the ECM: Financial performance (ROA)

Table 2 presents the cointegrating and the ECT results. PMG is the most efficient and preferred estimator. Financial stability and the corporate governance index have a cointegrating and negative relationship. In the long run, an increase in financial stability reduces the corporate governance index of the selected financial institutions. Furthermore, the results imply that when financial stability increases, corporate governance practice could be ineffective. The result is consistent with Kiemo et al. (2019) and Gaganis et al (2020) who found a significantly negative nexus between financial stability and corporate governance index. However, inconsistent with Mallin (2010) and Wahba (2015), who found a positive correlation between financial stability and corporate governance index. Mutuma (2024) argues that corporate governance should provide oversight to management to maximise the institution's financial stability. The same results were observed between the CAR and the corporate governance index. When the CAR is increased, the corporate governance index experiences a reduction in the long run and is significant at a 0.05 significance level. The result shows that an increase in the CAR, in the long run, widens the corporate governance index gap. According to Jensen and Meckling (1976), while agency theory aims to reduce agency costs, which may lead to a reduced freedom of agents, it may restrict managerial adaptability and initiatives. Furthermore, meanwhile corporate governance is important for financial institutions, the agency relationship between agents and principals may encourage excessive risk taking if agents prefer higher returns, which may significantly reduce CAR and financial stability. The association between ROA and corporate governance is insignificant.

**Table 2: Summary of the cointegrating results and the ECT: GOVINDEX**

| Variables | PMG D.GOVINDEX | MG D.GOVINDEX | DFE D.GOVINDEX |
|---|---|---|---|
| Long-run | | | |
| FINSTAB | -0.0219*** | 0.226 | 0.00377 |
| | (-5.21) | (0.48) | (0.21) |
| ROA | 0.00178 | 0.628 | 0.00964 |
| | (0.39) | (1.17) | (0.64) |
| CAR | -0.00536* | 0.793 | -0.00598 |
| | (-2.32) | (1.72) | (-0.89) |
| | | | |
| ECT | -0.665*** | -0.839*** | -0.464*** |
| | (-8.07) | (-9.19) | (-8.65) |
| Short run | | | |
| D.FINSTAB | -0.0248 | -0.0572 | -0.00718 |
| | (-0.48) | (-0.35) | (-0.99) |
| D.ROA | -0.00675 | -0.0485 | -0.00229 |
| | (-0.07) | (-0.34) | (-0.43) |
| D.CAR | -0.0169 | -0.322 | -0.00519* |
| | (-0.34) | (-1.42) | (-1.99) |
| _cons | 0.350** | -0.577 | 0.0245 |
| | (2.90) | (-0.74) | (0.23) |
| N | 273 | 273 | 273 |
| Hausman Test (MG & MPG) | 3.72 | 3.72 | - |
| Hausman Test (DFE & MPG) | 0.23 | - | 0.23 |

*** $p < 0.001$,** $p < 0.01$,* $p < 0.05$. Standard errors in parentheses. GOVINDEX (corporate governance proxies: BS, NED, INED, BR, BD, and TD), FINSTAB (financial stability), ROA (financial performance) and CAR (risk appetite). D. represents the difference operator.

**Table 3: Summary of the cointegrating results and the ECT: FINSTAB**

| Variables | PMG D.FINSTAB | MG D.FINSTAB | DFE D.FINSTAB |
|---|---|---|---|
| Long-run | | | |
| GOVINDEX | 0.0668 | 0.490 | 1.659 |
| | (1.08) | (0.20) | (1.86) |
| ROA | 0.256*** | 4.240* | -0.0214 |
| | (13.35) | (2.00) | (-0.21) |
| CAR | 0.244*** | 1.285 | 0.268*** |
| | (131.69) | (1.93) | (7.33) |
| | | | |
| ECT | -0.450*** | -1.177*** | -0.688*** |
| | (-4.71) | (-6.76) | (-9.64) |
| Short-run | | | |
| D.GOVINDEX | -3.405 | -2.262* | -1.707** |
| | (-1.36) | (-2.04) | (-2.67) |
| D.ROA | 3.307* | -0.615 | 0.0959 |
| | (2.55) | (-0.68) | (1.77) |
| D.CAR | -0.0773 | -2.159 | -0.0141 |
| | (-0.24) | (-1.24) | (-0.53) |
| _cons | 4.639** | 4.885 | 5.275*** |
| | (2.77) | (1.70) | (5.02) |
| N | 273 | 273 | 273 |
| Hausman Test (MG & MPG) | 0.99 | 0.99 | - |
| Hausman Test (DFE & MPG) | 49.60*** | - | 49.60*** |
| Hausman Test (MG & DFE) | - | 0.26 | 0.26 |

*** $p < 0.001$, ** $p < 0.01$,* $p < 0.05$. GOVINDEX (corporate governance proxies: BS, NED, INED, BR, BD, and TD), FINSTAB (financial stability), ROA (financial performance) and CAR (risk appetite). D. represents the difference operator.

Several studies, such as Joe-Duke (2011), Aziz (2021) and Aluoch (2023) found an insignificant link between ROA and corporate governance. Corporate governance measures are process-oriented, while financial performance exhibits asset utilisation and operational efficiency. Bhagat and Bolton (2008) assert that corporate governance may not directly impact financial performance in the short run. ECT is significant and negative under the more efficient estimator (PMG). The findings infer a cointegrating relationship among the variables under analysis, namely, corporate governance index, financial stability, ROA, and CAR, but more

so, -0.665 represents the speed of adjustment. Therefore, the speed of adjustments to equilibrium will be 66.5 percent per year, which suggests a moderately fast adjustment speed. While financial stability and CAR are statistically insignificant in the short run, their ECT indicates a significant existence of an essential long run relationship with corporate governance index. Table 2 presents the cointegrating and the ECT results. DFE is the more efficient and preferred estimator.

DFE is more efficient for financial stability; therefore, the results discussed are based on the DFE estimator. There is a cointegrating relationship between the CAR and financial stability. The relationship is positive and statistically significant. However, in the long run, an increase in the CAR will increase the financial stability of the financial institutions. Therefore, it contributes to a well-functioning and efficient financial system sector. Nguyen (2021) and Affes and Jarboui (2023) assert that an improvement in corporate governance will also increase the financial stability of financial institutions. However, the current study found that the associations between the corporate governance index and financial stability, as well as ROA and financial stability, are insignificant. These results are inconsistent with those of Antwi and Kwakye (2022), who found a positive and significant association between financial stability and ROA.

There is a cointegrating relationship among the variables under analysis: financial stability, corporate governance index, ROA and CAR. The cointegrating relationship in ECT is negative and significant. Therefore, the model is in disequilibrium. While CAR is statistically insignificant in the short run, the ECT indicates a significant existence of an essential long run relationship with financial stability. -0.688 represents the speed of adjustment, which implies the speed of adjustments to equilibrium will be 68.8 percent per year, which suggests a moderately fast adjustment speed. Meanwhile, the corporate governance index is significant in the short run, its impact fades in the long run due to loss of strategic agility.

Table 4 provides the results of the cointegrating relationship and the ECT. PMG is more efficient. Therefore, it is the preferred estimator, and the results are based on PMG.

### Table 4: Summary of the cointegrating results and the ECT: CAR.

| Variables | PMG D.CAR | MG D.CAR | DFE D.CAR |
|---|---|---|---|
| Long-run | | | |
| GOVINDEX | -0.0557 | -3.228 | -1.161 |
| | (-0.56) | (-1.66) | (-0.57) |
| FINSTAB | 0.230*** | 3.581** | 1.266*** |
| | (14.95) | (3.05) | (5.85) |
| ROA | 0.670*** | -0.255 | 0.894*** |
| | (7.49) | (-0.47) | (4.10) |
| | | | |
| ECT | -0.157* | -0.793* | -0.683*** |
| | (-2.61) | (-2.56) | (-12.25) |
| Short-run | | | |
| D.ROA | -0.851* | 0.254 | -0.321** |
| | (-2.25) | (0.58) | (-2.70) |
| D.GOVINDEX | -2.806 | -2.085 | -3.273* |
| | (-1.62) | (-1.23) | (-2.30) |
| D.FINSTAB | 2.205 | -1.033 | -0.0280 |
| | (1.67) | (-1.08) | (-0.17) |
| _cons | 14.45* | 5.676 | 5.684* |
| | (2.52) | (0.97) | (2.35) |
| N | 273 | 273 | |
| Hausman Test (MG & MPG) | 3.13 | 3.13 | - |
| Hausman Test (DFE & MPG) | 6.08 | - | 6.08 |

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$. GOVINDEX (corporate governance proxies: BS, NED, INED, BR, BD, and TD), FINSTAB (financial stability), ROA (financial performance) and CAR (risk appetite). D. represents the difference operator.

There is a cointegrating relationship between financial stability and the CAR. The long-run relationship is positive and significant at a 0.001 significance level. The higher the financial stability is, the higher the CAR. The results are consistent with the findings of Nguyen (2021) in Vietnamese financial institutions, where the relationship between the CAR and financial stability was positive and significant. When the study measured the relationship between the ROA and the CAR, the results showed a cointegrating relationship. The long-run relationship between the ROA and the CAR is positive and significant. The results imply that the higher the ROA, the higher the CAR for financial institutions. The CAR measures the financial institution's ability to meet its financial obligations by comparing its capital with its assets. The results of the current study are consistent with those of Shabani et al. (2019) and Benvenuto et al. (2021), who found a statistically significant and positive association between ROA and CAR. However, this is inconsistent with Setiawan and Irfani (2024), who found a negative link between ROA and CAR. The current study found an insignificant association between corporate governance and CAR.

The ECT is negative and significant. Therefore, the results confirm a cointegrating relationship between the variables (CAR, ROA, corporate governance index, and financial stability) under analysis, at a speed of adjustments to the equilibrium of 15.7 percent per year, which suggests a slow speed of adjustment. However, financial stability is insignificant in the short run, its impact increases

in the long run, indicating a significant level. While ROA is statistically significant in the short run, its ECT indicates a significant existence of an essential long run relationship with CAR. Table 5 provides the results of the cointegrating relationship and the ECT. PMG is more efficient. Therefore, it is the preferred estimator, and the results are based on PMG.

**Table 5: Summary of the cointegrating results and ECT: ROA**

| Variables | PMG D.ROA | MG D.ROA | DFE D.ROA |
|---|---|---|---|
| Long-run | | | |
| GOVINDEX | -0.00871 | -5.099 | -1.599 |
| | (-0.28) | (-1.63) | (-1.80) |
| FINSTAB | 0.0190*** | 2.895** | 0.0569 |
| | (3.42) | (2.99) | (0.47) |
| CAR | 0.0823*** | -0.597*** | 0.102* |
| | (15.02) | (-4.22) | (2.31) |
| ECT | -0.483*** | -0.959*** | -0.831*** |
| | (-4.75) | (-13.15) | (-12.99) |
| Short-run | | | |
| D.GOVINDEX | 0.617 | 1.309 | 1.649* |
| | (1.32) | (1.41) | (2.15) |
| D.FINSTAB | 2.495** | 0.344 | 0.0679 |
| | (2.60) | (1.11) | (0.77) |
| D.CAR | -0.108 | 0.259*** | -0.000803 |
| | (-1.32) | (3.84) | (-0.03) |
| _cons | 0.554* | -2.176 | -0.0142 |
| | (2.03) | (-1.35) | (-0.01) |
| N | 273 | 273 | 273 |
| Hausman Test (MG & MPG) | 0.81 | 0.81 | - |
| Hausman Test (DFE & MPG) | 6.35 | - | 6.35 |
| Hausman Test (DFE & MPG) | - | 0.56 | 0.56 |

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$. GOVINDEX (corporate governance proxies: BS, NED, INED, BR, BD, and TD), FINSTAB (financial stability), ROA (financial performance) and CAR (risk appetite). D. represents the difference operator.

The cointegrating relationship between ROA and the corporate governance index is insignificant. Ehikioya (2009) and Mollah et al. (2012) also found a negative and insignificant relationship between ROA and corporate governance index. The insignificant results suggest that strict corporate governance measures do not always affect the financial performance of financial institutions. There is a cointegrating relationship between financial stability and ROA. However, the cointegrating relationship is positive and significant. A percentage increase in the financial stability of the financial institutions increases the financial performance (ROA). This finding is consistent with Tan and Anchor (2016) and Antwi and Kwakye (2022), who found a significant and positive relationship between financial stability and ROA. However, the study found an insignificant association between corporate governance and ROA.

There is also a cointegrating relationship between the CAR and ROA. The relationship is both positive and significant. An increase in the financial institutions' CAR increases the ROA of the selected institutions. The result is consistent with those of Shabani *et al.* (2019), Mbaeri *et al.* (2021), and Benvenuto *et al.* (2021), who found positive and significant results. The ECT of the variables under analysis is negative and significant. Therefore, a cointegrating relationship exists between the variables under analysis: ROA, corporate governance index, financial stability, and CAR. The speed of adjustment to equilibrium will be 48.3 percent per year, which suggests a moderate speed of adjustment. However, CAR is insignificant in the short run, its impact increases in the long run, indicating a significant level. While financial stability is statistically significant in the short run, its ECT indicates a significant existence of an essential long run relationship with ROA.

## 4.2 Panel Cointegration and the ECM: Financial performance (ROE)
Table 6 provides the results of the cointegrating relationship and the ECT. DFE is more efficient. Therefore, it is the preferred estimator, and the results are based on DFE. This study found no cointegrating relationships between financial stability and the corporate governance index, ROE and the corporate governance index, CAR and the corporate governance index. Ajanthan et al. (2013), Elbahar (2016), and Bawaneh (2020) found an insignificant link between corporate governance and financial performance. However, the results are contrary to the agency theory prediction that corporate governance improves financial stability and performance (Jensen and Meckling, 1976). The cointegrating relationship between financial stability and the corporate governance index is positive but statistically insignificant in the long run. A cointegrating relationship exists among the variables under analysis, namely, corporate governance index, financial stability, ROE, and CAR. The model is in disequilibrium, and the speed of adjustments to equilibrium is 47.4 percent per year, which suggests a moderate speed of adjustment. However, CAR, ROE and financial stability are all insignificant in the short run, their impact also does not significantly increase in the long run. Table 7 provides the results of the cointegrating relationship and the ECT. DFE is more efficient. Therefore, it is the preferred estimator, and the results are based on DFE.

**Table 6: Summary of the cointegrating results and the ECT: GOVINDEX**

| Variables | PMG D.GOVINDEX | MG D.GOVINDEX | DFE D.GOVINDEX |
|---|---|---|---|
| Long-run | | | |
| FINSTAB | -0.0191*** | 0.288 | 0.00315 |
| | (-5.16) | (0.68) | (0.18) |
| ROE | -0.000184 | -0.0237 | -0.00290 |
| | (-0.08) | (-0.96) | (-0.38) |
| CAR | -0.00521* | 0.511 | -0.00427 |
| | (-2.56) | (1.43) | (-0.67) |
| | | | |
| ECT | -0.669*** | -0.842*** | -0.474*** |
| | (-8.25) | (-9.35) | (-8.85) |
| Short-run | | | |
| D.FINSTAB | 0.0229 | -0.0668 | -0.00724 |
| | (0.70) | (-0.42) | (-1.00) |
| D.ROE | -0.00812 | -0.00645 | 0.00107 |
| | (-1.08) | (-0.67) | (0.36) |
| D.CAR | -0.0246 | -0.190 | -0.00529* |
| | (-0.69) | (-1.18) | (-2.04) |
| _cons | 0.323** | -0.139 | 0.0468 |
| | (2.72) | (-0.21) | (0.38) |
| N | 273 | 273 | 273 |
| Hausman Test (MG & MPG) | 4.54 | 4.54 | - |
| Hausman Test (DFE & MPG) | 38.85*** | - | 38.85*** |
| Hausman Test (MG & DFE) | - | 0.25 | 0.25 |

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$. GOVINDEX (corporate governance proxies: BS, NED, INED, BR, BD, and TD), FINSTAB (financial stability), ROE (financial performance) and CAR (risk appetite). D. represents the difference operator.

**Table 7: Summary of the cointegrating results and ECT: FINSTAB**

| Variables | PMG D.FINSTAB | MG D.FINSTAB | DFE D.FINSTAB |
|---|---|---|---|
| Long-run | | | |
| GOVINDEX | 0.494** | -2.051 | 1.400 |
| | (2.82) | (-0.60) | (1.59) |
| ROE | 0.197*** | 0.0609 | -0.0714 |
| | (11.73) | (0.29) | (-1.33) |
| CAR | 0.269*** | 1.474 | 0.265*** |
| | (312.01) | (1.94) | (7.99) |
| | | | |
| ECT | -0.462*** | -1.142*** | -0.704*** |
| | (-5.50) | (-14.63) | (-9.82) |
| Short-run | | | |
| D.GOVINDEX | -3.852 | -3.204 | -1.589* |
| | (-1.16) | (-1.71) | (-2.48) |
| D.ROE | 0.0383 | 0.00127 | 0.0450 |
| | (0.86) | (0.02) | (1.49) |
| D.CAR | -0.0322 | -2.095 | -0.00946 |
| | (-0.12) | (-1.62) | (-0.35) |
| _cons | 3.262* | 6.151 | 6.213*** |
| | (2.09) | (1.34) | (5.14) |
| N | 273 | 273 | 273 |
| Hausman Test (MG & MPG) | 1.63 | 1.63 | - |
| Hausman Test (DFE & MPG) | 40.00*** | - | 40.00*** |
| Hausman Test (MG & DFE) | - | 1.23 | 1.23 |

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$. GOVINDEX (corporate governance proxies: BS, NED, INED, BR, BD, and TD), FINSTAB (financial stability), ROE (financial performance) and CAR (risk appetite). D. represents the difference operator.

The cointegrating relationship between financial stability and corporate governance is positive and insignificant. Similar results were found by Haribowo et al. (2021), who found no significant impact of corporate governance on financial stability. A cointegrating relationship exists between the CAR and the financial stability of financial institutions. The long-run relationship is positive and significant at a 0.001 significance level. The results are consistent with those of Nguyen (2021), who found a positive and significant relationship between the CAR and financial stability. The result implies that an increase in the CAR will increase the financial stability of the financial institutions. However, the current study found an insignificant association between the corporate governance index and financial stability and between ROE and financial stability. The ECT is negative but highly significant at 0.001. Therefore, a cointegrating relationship between the variables exists, namely, financial stability, corporate governance index, ROE, and CAR under analysis, with -0.704 representing the speed of adjustment. Therefore, the speed of adjustments to equilibrium will be 70.4 percent per year, which suggests a moderately fast speed of adjustment. However, CAR is insignificant in the short run, its impact increases in the long run, reaching a significant level.

Table 8 provides the results of the cointegrating relationship and the ECT. DFE is more efficient. Therefore, it is the preferred estimator, and the results are based on DFE.

**Table 8: Summary of the cointegrating results and ECT: CAR.**

| Variables | PMG D.CAR | MG D.CAR | DFE D.CAR |
|---|---|---|---|
| Long-run | | | |
| GOVINDEX | 0.663 | -2.785 | -1.574 |
| | (0.92) | (-1.43) | (-0.72) |
| FINSTAB | 0.736*** | 3.337** | 1.482*** |
| | (6.01) | (2.93) | (6.63) |
| ROE | 0.508*** | 0.0278 | 0.331* |
| | (7.04) | (0.17) | (2.51) |
| | | | |
| ECT | -0.310* | -1.739** | -0.648*** |
| | (-2.20) | (-2.72) | (-11.59) |
| Short-run | | | |
| D.ROE | -0.184* | -0.0378 | -0.156* |
| | (-2.17) | (-0.48) | (-2.31) |
| D.GOVINDEX | -2.733 | -1.879* | -2.849* |
| | (-1.62) | (-2.01) | (-1.97) |
| D.FINSTAB | 1.974 | -0.746 | -0.0655 |
| | (1.69) | (-1.00) | (-0.39) |
| _cons | 9.189 | 7.443 | 2.134 |
| | (1.67) | (1.18) | (0.75) |
| N | 273 | 273 | 273 |
| Hausman Test (MG & MPG) | 3.27 | 3.27 | - |
| Hausman Test (DFE & MPG) | | | |
| Hausman Test (MG & DFE) | 3.13 | - | 3.13 |
| | - | 55.75*** | 55.75*** |

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$. GOVINDEX (corporate governance proxies: BS, NED, INED, BR, BD, and TD), FINSTAB (financial stability), ROE (financial performance) and CAR (risk appetite). D. represents the difference operator.

There is no cointegrating relationship between CAR and corporate governance. Benvenuto et al. (2021) found similar results where there was no significant cointegrating relationship between CAR and corporate governance. There is a cointegrating relationship between financial stability and the CAR. The long-run relationship is positively significant at 0.001. In the long run, financial stability will increase financial institutions' CAR. This finding is consistent with the results of Nguyen (2021) and Benvenuto *et al*. (2021), who found a significant and positive relationship between financial stability and the CAR. However, the current study found an insignificant association between corporate governance and CAR.

The relationship between ROE and the CAR is also significant and positive. Therefore, there is a cointegrating relationship between these variables. The higher the ROE, the higher the CAR of the selected financial institutions. The results are consistent with those of Angahar *et al*. (2019) and Shabani *et al*. (2019), who found a positive and significant relationship between ROE and the CAR. Financial institutions use the CAR to assess the sufficiency of their capital holdings in light of their exposures. The ECT is negative and significant. Therefore, a cointegrating relationship exists between the variables: CAR, corporate governance index, financial stability, and ROE. Under this analysis, the speed of adjustments to equilibrium will be 31 percent per year, which suggests a moderate adjustment speed. While financial stability is statistically insignificant in the short run, the ECT indicates a significant existence of an essential long run relationship with CAR. ROE is significant in the short run, its impact also increases in the long run due to increase in strategic agility. However, corporate governance is only significant in the short run, and its impact fades in the long run.

Table 9 provides the results of the cointegrating relationship and the ECT. PMG is more efficient. Therefore, it is the preferred estimator, and the results are based on PMG. The long-run relationship between the corporate governance index and ROE is insignificant. The results imply that corporate governance measures have no significant long-run relationship with ROE. Furthermore, corporate governance loses relevance over time. Furthermore, the relationship between financial stability and ROE is

also insignificant. Moreover, the relationship between the CAR and ROE is insignificant. Under the preferred PMG estimator, the ECT is negative and statistically significant at 0.001. Therefore, a cointegrating relationship between the variables exists: ROE, corporate governance index, financial stability, and CAR under analysis. The speed of adjustment to equilibrium will be 57.7 percent per year, which suggests a moderate adjustment speed. While financial stability is statistically significant in the short run, its relevance losses in the long run.

**Table 9: Summary of the cointegrating results and ECT: ROE.**

| Variable | PMG D.ROE | MG D.ROE | DFE D.ROE |
|---|---|---|---|
| Long-run | | | |
| GOVINDEX | -0.138 | 62.65 | -2.759 |
| | (-0.21) | (0.93) | (-1.61) |
| FINSTAB | 0.128 | 6.283 | 0.000331 |
| | (1.14) | (1.40) | (0.00) |
| CAR | -0.0433 | -4.543 | 0.00410 |
| | (-1.40) | (-1.15) | (0.05) |
| | | | |
| ECT | -0.577*** | -0.969*** | -0.783*** |
| | (-6.37) | (-11.76) | (-12.19) |
| Short-run | | | |
| D.GOVINDEX | 2.579 | 6.133* | 2.013 |
| | (1.44) | (2.03) | (1.45) |
| D.FINSTAB | 6.196* | -2.267 | -0.0248 |
| | (2.18) | (-0.98) | (-0.15) |
| D.CAR | 1.283 | 6.872 | 0.0505 |
| | (0.57) | (1.55) | (0.87) |
| _cons | 7.903*** | 7.122 | 12.18*** |
| | (5.84) | (1.24) | (4.62) |
| N | 273 | 273 | 273 |
| Hausman Test (MG & MPG) | 0.73 | 0.73 | - |
| Hausman Test (DFE & MPG) | 0.60 | - | 0.60 |
| Hausman Test (MG & DFE) | | | |
| | - | 4.17 | 4.17 |

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$. GOVINDEX (corporate governance proxies: BS, NED, INED, BR, BD, and TD), FINSTAB (financial stability), ROE (financial performance) and CAR (risk appetite). D. represents the difference operator.

As a result of the entire test, in which the corporate governance index was the dependent variable, the ECT, measuring the speed of adjustments for long-run equilibrium, is significant and negative. ECT must be significant and negative to correct the short-run divergence to its long-run equilibrium (Gujarati, 2021). None of the ECTs in this study were positive, which implies that the time series converges towards the long-term equilibrium. The results satisfied the PMG and DFE conditions of the long-run relationships. The negative and significant coefficients of ECT were between 0 and -1.

Using corporate governance proxies as dependent variables, the study discussed the cointegration relationships. Therefore, the current study reported cointegrating relationships between the chosen independent variables.

## 5. Conclusions

The study was limited to South African financial institutions registered under the FSCA and the Bureau Van Dijk Orbis Bank, incorporating data from 2007 to 2020. The study included other financial variables (dimensions), namely, financial stability, risk appetite, and financial performance, in investigating the cointegrating association between financial performance and corporate governance. We employed the principal component analysis method to develop a composite index to proxy corporate governance instead of only using the individual corporate governance proxies: board diversity, remuneration, composition, and size. Therefore, the corporate governance index of the current study was necessary to capture and reflect the corporate governance differences in the sample of selected financial institutions. Moreover, incorporating the corporate governance index into the study further emphasised its importance in financial institutions.

When financial stability was regressed as the dependent variable, we concluded that financial stability in the selected financial institutions had cointegrating relationships with the corporate governance index, CAR, and ROA when the financial performance measure was used as ROA. Furthermore, financial stability cointegrated with the CAR when the financial performance measure was used as ROE. When the CAR was regressed as the dependent variable employing ROA to measure financial performance, we found cointegrating relationships between the CAR and financial stability, and between the CAR and ROA.

However, capital adequacy (CAR) had a cointegrating relationship with financial stability when the financial performance measure was ROE. When the ROA was regressed as the dependent variable, we found a cointegrating relationship between the ROA and the CAR. When the ROE was regressed as the dependent variable, we found a cointegrating relationship between the ROE and CAR. The presence of a cointegration relationship means that there is a long-term equilibrium between the variables.

Similar to Min *et al*. (2015), we concur that governance systems for financial institutions should have been created before the global financial crisis of 2007. However, regulators learned great lessons, and notable progress and improvements have been made since

then, although there are still some loopholes. We thus recommend that greater attention be paid to the enterprise risk management of financial institutions to enable them to identify and set their risk appetite thresholds, which impact the profitability of their operations. Regulators must also ensure that they continuously adjust the capital adequacy ratios of financial institutions, in line with their respective risk profiles, to avoid bank runs and failures.

Future research could consider applying Tobin's Q as a finance measure when gauged against individual corporate governance variables, as was the case in the study by Park and Byun (2022). Such results could then be compared to determine whether it matters that one study used a composite index, while others used individual variables.

# References

Affes, W. and Jarboui, A. (2023). The impact of corporate governance on financial performance: a cross-sector study. *International Journal of Disclosure and Governance*, 20(4): 374-394.

Ajanthan, A., Balaputhiran, S. and Nimalathashan, B. (2013). Corporate governance and banking performance: A comparative study between private ad state banking sector in Sri Lanka. *European Journal of Business and Management*, 5(20): 92-100.

Angahar, P.A., Tsokar, E.E. and Agbo, A. (2019). Effect of capital structure on corporate performance of listed consumer goods firms in Nigeria, *Journal of Accounting, Finance and Management Discovery,* 2(2): 1-119.

Animasaun, R.O., Omotunwase, O.M., Babayanju, A.G.A. and Bamgboye, A.A. (2025). Effect of credit risk management on financial performance of listed deposit money banks in Nigeria. *International Journal of Research in Social Science and Humanities,* 6(1): 1-12.

Antwi, F. and Kwakye, M. (2022). Modelling the effect of bank performance on financial stability: Fresh evidence from Africa, *Research in Business & Social Science,* 11(7): 2147-4478.

Aluoch, M. O. (2023). Corporate governance and performance of commercial banks listed at the Nairobi Securities Exchange, Kenya. *European Scientific Journal, ESJ*, 19(10), 194.

Apergis, N. and Payne, J.E. (2010). Renewable energy consumption and economic growth: Evidence from a panel of OECD countries, Energy policy 38(1): 656–660.

Aziz, S. and Abbas, U. (2021). Corporate governance, innovation and corporate performance: A study of pharmaceutical industry of Pakistan. *Asian Journal of Economics and Finance*, 3(4), 459-471.

Baltagi, B. (2008). *Econometric analysis of panel data*. Hoboken, NJ: Wiley.

Benvenuto, M., Avram, R.L., Avram, A. and Viola, C. (2021). Assessing the impact of corporate governance index on financial performance in the Romanian and Italian banking systems. *Sustainability*, 13(10): 1-16.

Biresaw, T.M. and Sibindi, A.B. (2025). A Structured Approach to Enterprise Risk Management in Ethiopian Commercial Banks. *Risk Management Magazine*, 20(1): 57-77. https://www.aifirm.it/wp-content/uploads/2025/04/RMM-2025-01-Excerpt-4.pdf

Bhagat, S. and Bolton, B. (2008). Corporate governance and firm performance. *Journal of corporate finance*, 14(3), 257-273.

Brealey, R. A., Myers, S. C., Allen, F. and Edmans, A. (2022). *Principles of corporate finance.* 14th edition. New York: McGraw-Hill Irwin.

Croissant, Y. and Millo, G., 2019. *Panel data econometrics with R*. Wiley.

Donaldson, L. and Davis, J. H. (1991). Stewardship theory or agency theory: CEO governance and shareholder returns. *Australian Journal of Management*, 16(1): 49–64.

Efunniyi, C.P., Abhulimen, A.O., Obiki-Osafiele, A.N., Osundare, O.S., Agu, E.E. and Adeniran, I.A. (2024). Strengthening corporate governance and financial compliance: Enhancing accountability and transparency. *Finance & Accounting Research Journal*, 6(8): 1597-1616.

Ehikioya, B. I. (2009). Corporate governance structure and firm performance in developing economies: evidence from Nigeria. *Corporate Governance: The international journal of business in society*, 9(3), 231-243.

Elbahar, E. (2016). *Corporate governance, risk management, and bank performance in the GCC banking sector.* Thesis. University of Plymouth. Available at: https://doi.org/10.24382/1002

Engle, R.F. and Yoo, B.S. (1987). Forecasting and testing in co-integrated systems. *Journal of Econometrics,* 35(1): 143-159.

Fujiki, H., Hsiao, C. and Shen, Y. (2002). Is there a stable money demand function under the low interest rate policy? A panel data analysis, *Monetary & Economic Studies*, 20(2), 1–23.

Gaganis, C., Lozano-Vivas, A., Papadimitri, P. and Pasiouras, F. (2020). Macroprudential policies, corporate governance and bank risk: Cross-country evidence. *Journal of Economic Behavior & Organization*, 169(1): 126-142. https://doi.org/10.1016/j.jebo.2019.11.004

Greenacre, M., Groenen, P.J., Hastie, T., D'Enza, A.I., Markos, A. and Tuzhilina, E. (2022). Principal component analysis. *Nature Reviews Methods Primers*, 2(1), 1-21.

Gujarati, D.N. (2021). *Essentials of econometrics*. 5th ed. SAGE Publications.

Gwaison, P.D. and Maimakp, L.N. (2021). Effects of corporate governance on financial performance of commercial banks in Nigeria, *International Journal of Financial Research,* 2(1): 13-23.

Haribowo, I., Putri, Z. E. and Yulianti, Y. (2021). Corporate Governance and Financial Stability of Islamic Banks in Asia. *The Journal of Asian Finance, Economics and Business*, 8(12): 353-361.

Hoffman, D.L. and Rasche, R.H. (1996). Assessing forecast performance in a cointegrated system. *Journal of Applied Econometrics,* 11(5):495-517.

Hsiao, C. (2022). *Analysis of panel data*. Cambridge University Press.

Hsiao, C., Mountain, D.C. and Illman, K.H. (1995). A Bayesian integration of end-use metering and conditional-demand analysis, Journal of Business & Economic Statistics 13(3), 315–326.

Hunjra, A.I., Jebabli, I., Thrikawala, S.S., Alawi, S.M. and Mehmood, R., 2024. How do corporate governance and corporate social responsibility affect credit risk?. *Research in International Business and Finance*, *67*: 1-18. https://doi.org/10.1016/j.ribaf.2023.102139

Jensen, M. C. and Meckling, W. H. (1976). Theory of the Firm: Managerial Behaviour, Agency Costs, and Ownership Structure. *Journal of Financial Economics* 3(4): 305–50. https://doi.org/10.1016/0304-405X(76)90026-X

Joe-Duke, I. I. and Kankpang, K. A. (2011). Linking corporate governance with organizational performance: New insights and evidence from Nigeria. *Global Journal of Management and Business Research*, 11(12), 46-58.

Jouirou, M. and Jouini, F. (2022). Corporate governance mechanisms and banking performance. *Global Business and Economics Review*, *27*(4): 475-491.

Kajumbula, R. and Makoni, P.L. (2024). CEO power and bank risk nexus: Evidence from commercial banks in Uganda. *Risk Management Magazine*, 19(2): 42-53. http://dx.doi.org/10.47473/2020rmm0143.

Karpoff, J.M. (2021). On a stakeholder model of corporate governance. *Financial Management*, *50*(2): 321-343.

Khoza, F., Makina, D. and Makoni, P. L. (2024). Key determinants of corporate governance in financial institutions: evidence from South Africa; Risks 12 (6):90. https://doi.org/10.3390/risks12060090

Kiemo, S.M., Olweny, T.O., Muturi, W.M. and Mwangi, L.W. (2019). Bank-specific determinants of commercial banks financial stability in Kenya. *Journal of Applied finance and banking*, 9(1): 119-145.

Kiptoo, I. K., Kariuki, S. N., Ocharo, K. N. and Ntim, C. G. (2021). Corporate governance and financial performance of insurance firms in Kenya. Cogent Business & Management, 8(1): 1-17. https://doi.org/10.1080/23311975.2021.1938350

Lingwati, E. and Mamabolo, A. (2023). Emerging market entrepreneurs' narratives on managing business ethical misconduct. *Acta Commercii-Independent Research Journal in the Management Sciences*, *23*(1): 1-12.

Mahmudi, B. (2024). Corporate governance mechanisms and financial performance: A systematic literature review in emerging markets. *Management Studies and Business Journal,* 1(3): 270-285.

Mallin, C.A. (2010). Corporate governance and risk: A study of board structure and process. *Journal of Financial Regulation and Compliance*, 18(2): 174-191.

Mbaeri, M.N., Uwalake, U. and Gimba, J.T. (2021). Capital adequacy ratio and financial performance of listed commercial banks in Nigeria. *Journal of Economics and Allied Research*, *6*(3): 81-88.

Min, B. S., Cashel-Cordo, P. and Rhim, J. C. (2015). Corporate leverage strategy in an emerging market. *Global Business & Finance Review*, 20(2), 35-48.

Mollah, S., Al Farooque, O. and Karim, W. (2012). Ownership structure, corporate governance and firm performance: Evidence from an African emerging market. *Studies in Economics and finance*, 29(4), 301-319.

Msomi, T. S. and Nzama, S (2023). Analyzing firm-specific factors affecting the financial performance of insurance companies in South Africa. *Insurance Markets and Companies*, 14(1): 8-21. doi:10.21511/ins.14(1).2023.02

Musa, H. (2020). Corporate governance and financial performance of Nigeria listed banks. *Journal of Advanced Research in Dynamical & Control Systems*, 12(1): 5-10.

Mutuma, F. D. K. W. (2024). Emerging trends in board leadership and corporate governance. *ICS Governance Journal*, *1*(3): 23-37.

Muzata, T. and Marozva, G. (2023). The Nexus between Executive Compensation and Firm Performance: Does Governance and Inequality Matter? *Global Business & Finance Review*, 28(5): 31-50

Nguyen, M.S. (2021). Capital adequacy ratio and a bank's financial stability in Vietnam. *Banks and Bank Systems*, *16*(4): 61.

Nizam, K. and Liaqat, O. (2022). Corporate Governance and firm performance: empirical evidence from Pakistan banking sector. *International Journal of Business Diplomacy and Economy*, *1*(2): 18-28.

Nyaloti, L.A. (2024). Effect of stock market crashes in international markets on the global economy: A case study of US stock market crashes. *International Journal of Research publication and Reviews*, 5(8): 4145-4150.

Oladipupo, O.E. and Kelvin, A.O. (2024). Corporate governance and manufacturing firms' financial performance in Nigeria. *Asian Journal of Economics, Business and Accounting,* 24(11): 471-490.

Owusu, C.K. and Garr, D.K. (2024). Corporate governance dynamics and financial performance: Analysis of listed commercial banks in the Ghanaian context. *International Journal of Business, Management, and Economics,* 5(2): 114-133.

Park, K. H. and Byun, J. (2022). Board diversity, IPO underpricing, and firm value: evidence from Korea. *Global Business & Finance Review*, 27(1), 65-82

Pesaran, M. H., Shin, Y., and Smith, R. P. (1999). Pooled mean group estimation of dynamic heterogeneous panels. Journal of the American Statistical Association, 94(446), 621-634.

Phillips, P. C. (1991). Optimal inference in cointegrated systems. *Econometrica: Journal of the Econometric Society*, 1(1): 283-306.

Schroeck, G. (2002). *Risk management and value creation in financial institutions*. G. Schroeck, (Ed.). U.S.: Wiley.

Setiawan, F. and Irfani, A. F. (2024). Predicting capital adequacy ratio of Islamic rural banks based on FDR, NPF, ROA, and BOPO. In *Proceeding International Conference on Law, Economy, Social and Sharia (ICLESS),* 2: 748-771.

Shabani, H., Morina, F. and Misiri, V. (2019). The effect of capital adequacy on returns of assets of commercial banks in Kosovo. *European Journal of Sustainable Development*, *8*(2): 201-201.

Simanjuntak, B. and Alfredo, J. (2024). The impact of corporate governance on financial performance Evidence from Indonesia. *International Journal of Strategic Accounting and Business Management,* 1(1): 1-5.

Swedan, Z. M. and Ahmed, S. (2019). The impact of corporate governance implementation on firm's performance in the Gulf countries (GCC) based on stewardship assumptions. *International Journal of Business Society*, *2*(11): 42-47.

Talatu, H. F. (2024). Effect of corporate governance on financial performance of quoted healthcare firms in Nigeria. *ANUK College of Private Sector Accounting Journal,* 1(2): 69-78.

Tan, A. Y. and Anchor, J. R. (2016). Stability and profitability in the Chinese banking industry: evidence from an auto-regressive-distributed linear specification. *Investment Management and Financial Innovations*, *13*(4): 120-128.

Temba, G. I., Kasoga, P. S. and Keregero, C. M. (2023). Corporate governance and financial performance: Evidence from commercial banks in Tanzania. *Cogent Economics & Finance*, 11(2). https://doi.org/10.1080/23322039.2023.2247162

Tricker, B. and Tricker, R. I. (2015). *Corporate governance: Principles, policies, and practices*. Oxford University Press, USA.

Usendok, I. G. (2022). Corporate governance and organizational performance: A study of selected banks in Nigeria. *International Journal of Business and Management Review*, *10*(4): 59-74.

Wahba, H. (2015). The joint effect of board characteristics on financial performance: Empirical evidence from Egypt. *Review of Accounting and Finance*, *14*(1): 20-40.

# Ratings and Banking Regulation: a Shift from Productive (Basel II), to Contradictory (EBA-GL LOM and Supervisory Practices), to Dangerous (Basel 3+ and CRR3)

**Giacomo De Laurentis (Bocconi University, Italy)**

## Abstract

External ratings (agency ratings) and internal ratings were already in use before Basel II, as powerful management tools. Basel II "adopts" them (the former in the Standard approach, the latter in the IRB approaches) in a productive way: internal ratings represent the final summary of the creditworthiness assessment of debtors/transactions, to be used both in lending/review processes and in those aimed at regulatory capital adequacy measures, to be developed with a medium to long-term target horizon, leaving banks with the discretion to choose the assignment methods, as long as all relevant information (including qualitative and forward-looking data) is included in the judgment.

Instead, in the EBA-GL LOM and supervisory practices, internal ratings are predominantly the result of statistical tools in which the role of behavioural information is much broader compared to that of qualitative, strategic, and prospective information; they target short-term forecasting horizons, approaching early warning systems used in ongoing monitoring; and they are often definitively approved by bank officers different from those who underwrite the loans.

Basel 3+ confirms the Basel II framework, ignores the contradictions in the EBA-GL LOM and supervisory practices, and generates the false perception that internal ratings are no longer essential. In addition, CRR3 emphasizes the use of external ratings also in the SMEs segment and encourages the creation of new rating agencies (including those linked to central banks), thus pushing toward a commodity-oriented logic in bank-SMEs relationships (in contrast with the EBA-GL LOM).

One may wonder whether the evolution of the role and content of ratings in regulations and supervisory practices is the result of a deliberate, conscious shift in approach or if it is simply the outcome of the occasional predominance of differing positions, without any central authority, even a "Czar of ratings."[1]

## 1. Credit Rating Before Basel II: A Powerful Tool for Bank Management and Competitiveness.

In markets where banks are price setters, the use of ratings allows differentiation in the economic conditions required, lowering rates for better clients and increasing them for worse clients. This way, financing is expanded to less risky debtors/transactions and is left room to other banks to finance borrowers/transactions that are too poorly remunerative relative to risk (and to the future credit losses that will impact financial statements and deplete the lender's capital). In markets where the bank is price taker and must apply the "market rate," ratings allow the evaluation of the adequacy of returns relative to risk, meaning risky loans above the average (whose returns do not adequately cover the risk) are rejected, and less risky loans (on which the "margin after losses" is wide) are accepted. Essentially, ratings are a powerful competitive tool that causes significant damage to banks that do not use them, due to the phenomenon of "adverse selection" outlined above.

For these reasons, external ratings have been developed (by many of the current international rating agencies) since the early 20th century for the benefit of all investors, and internal ratings have been developed by the most advanced banks since the mid-1990s. These institutions began passing on the damages from adverse selection to other banks. The link between price and credit risk has steadily strengthened (De Laurentis, Maino, Molteni, 2010, § 7.4).

## 2. The Productive Relationship Between Regulation and Ratings Established by Basel II.

The reasons why the relationship between Basel II regulation (BIS, 2004) and ratings is extremely productive are three, with some corollaries.

**First reason: Incentives to improve bank risk management tools.**

To prevent the regulation itself from inducing adverse selection (by requiring, as Basel I did, the same capital absorption for private debtors and businesses of any risk level), Basel II welcomes external ratings in the calculation of First Pillar capital requirements for "Standard Banks." Additionally, by providing lower capital requirements for "IRB Banks," the regulation strongly incentivized the development and use of internal ratings. Finally, an additional incentive comes from Basel II's Second Pillar, as banks that do not have an essential tool for modern credit management, such as internal ratings, can face additional capital requirements defined on a bank-by-bank basis by the competent supervisory authority.

**Second reason: Internal ratings recognized as competitive management tools, the final synthesis of creditworthiness evaluations, and capable of incorporating all relevant information.**

The second reason why the relationship between Basel II and ratings is extremely productive is that the regulation emphasizes the role of internal ratings, developed autonomously by individual banks, and differing from bank to bank. This creates a significant burden for supervisory authorities to validate each rating system individually but has the great advantage of recognizing these systems as competitive tools, which banks continuously seek to improve to make more timely and accurate credit decisions than their competitors, even for the same debtors. The fundamental philosophy of Basel II is, in fact, to ensure that the systems used for loans underwriting decisions are the same systems used to calculate the capital requirements resulting from those decisions. Ratings are seen as dual-use tools: this is why the "use test" is required to validate systems for IRB approaches (par. 444 of Basel II (BIS, 2004) [(1)].

Thus, Basel II rejected the alternative hypothesis of introducing a "regulatory rating system." This would have quickly led to a divergence between ratings developed for management purposes to "beat" other banks and regulatory ratings used to determine capital

---

[1] A previous version of this article was published in Italian in issue no. 4 of 2025 of Bancaria, the journal of the Italian Banking Association.

requirements. Therefore, under Basel II, the rating is a competitive tool that allows banks that acquire more accurate, forward-looking client information and have more effective processes to synthesize it into ratings, to make better credit decisions.

A first corollary of all this is that the rating, in Basel II, represents the "final synthesis of creditworthiness assessment," on which the decision to assume a certain credit risk is based and to which the associated regulatory capital requirements are aligned. Therefore, the rating must include all the information considered in the lending decision regarding the debtor's/transaction's risk profile; evidently, this includes qualitative information and data related to the debtor's economic and financial prospects. Basel II specifically and explicitly confirms this corollary: in paragraph 417, regarding the possibility of using models for assigning ratings, it states: "However, there must be adequate evaluation and verification by those responsible to ensure that all relevant and pertinent information, including those outside the scope of the model, is considered and that it is used correctly." The Guide to Internal Models published by the ECB in October 2019, in par. 4.1.3 at point 64 letter a, confirms that "all relevant information should be included in the rating/grade/pool assignment process."

A second corollary is that banks must be free to choose how to assign ratings. In fact, Basel II is completely neutral regarding rating assignment methods; the regulation sets very stringent requirements only "downstream" of the "rating assignment processes" (risk differentiation), that is, for the "calibration processes" (risk quantification) of ratings.

Par. 417 of Basel II is particularly expressive: "The requirements set forth in this section apply to statistical models and other automatic methods for assigning ratings or estimating PD, LGD, and EAD. Credit scoring models and other automatic rating procedures generally use only a subset of available information. While they may sometimes avoid some of the errors typical of systems where subjective judgment plays an important role, the mechanical use of limited information is also a source of error. The models and procedures mentioned are permissible as a primary or partial basis for rating assignment and can contribute to estimating risk characteristics."

Thus, according to Basel II, models are not prohibited, but they are not mandatory either, nor is it obligatory to rely on them in a mechanistic way, limiting the role of overrides. Similarly, in European regulation (EU CRR, 2013, including Article 174).

**Third reason: Internal ratings as Tools for Granting/review processes, not Ongoing monitoring processes.**

The third reason why the relationship between Basel II and ratings is extremely productive concerns the exemplary clarity in the regulation regarding the role of ratings in credit processes.

The provision in paragraph 425 of Basel II "ratings must be updated at least annually" clearly indicates that ratings are considered tools to be used in granting/review processes, not in credit ongoing monitoring processes. Paragraph 414 of Basel II states: "Although the time horizon in the estimation of PD is one year, it is expected that banks use a longer time horizon when assigning ratings". Similarly, the ECB's Guide to Internal Models further specifies (ECB, 2019), in paragraph 4.1.3 "Grade assignment dynamics," at point 64 states: "Although the time horizon used in PD estimation is one year, it is the ECB's understanding that the rating/grade/pool assignment process should also adequately anticipate and reflect risk over a longer time horizon and take into account plausible changes in economic conditions. In order to achieve this objective: a) all relevant information should be included in the rating/grade/pool assignment process, giving an appropriate balance between drivers that are predictive only over a short time horizon and drivers that are predictive over a longer time horizon; b) a horizon of two to three years is considered to be appropriate for most portfolios."

It is clear that the time horizon to target when assigning ratings (and, thus, constructing statistical-based models if such methodologies are chosen) must be longer than one year, while the time horizon for calculating capital requirements to which the PD must refer (i.e., the probability of default associated with the ratings, once they have been assigned) is one year: this PD is the input to the risk-weighting functions that result in the RW (risk weight) of the bank's assets and determine, in turn, the regulatory capital required by regulation.

In essence, the regulation clearly distinguishes the time horizon to be used in the "quantification phase" of the ratings ("one year") from the time horizon to be targeted in the "assignment phase" of the ratings ("longer"). The latter meaning a) from the model development phase, when statistical-based rating models are built; b) in daily operations, when judgmental analyses are conducted by analysts to assign ratings.

**Contradictions in Basel II Regulation, Leading to Future Problems.**

Basel II does not lack contradictions regarding the construction methods and properties of ratings, particularly regarding the requirement that, on one hand, they should be forward-looking over non-short time horizons and stable over time (good risk differentiation), and, on the other hand, they should produce PD estimates that are close to actual default rates realized in specific future periods (good calibration, or quantification). These two properties of ratings can only be achieved with ratings of different natures.

The first is achieved with ratings defined as "Through the Cycle" (TTC), while the second is achieved with ratings called "Point in Time" (PIT). These two "rating philosophies" (as defined by the Basel Committee already in BIS, 2005; see also EBA, 2017, in the section titled "Rating philosophy") imply different rating system construction logics, both in the assignment and quantification processes, producing ratings with different characteristics. TTC ratings are stable over time because they assign judgments by looking at the debtor's prospects across an entire economic and sectoral cycle, particularly focusing on the debtor's reliability during the bottom of the cycle. Therefore, debtors tend not to change their rating class when the economy or sector improves or worsens. Moreover, PDs are associated with rating classes based on long-term historical default rate evidence. The consequence is that when comparing the PDs associated with a rating class to the actual default rates realized in subsequent periods for that class, it is discovered that the actual default rates are highly variable compared to the expectations incorporated into the estimated PD for that class. Thus, the rating has a very loose relationship with actual default rates. Every aspect is "reversed" in the PIT logic[2].

Regulation requires ratings to be assigned and quantified using TTC logic but also to exhibit the good calibration/quantification typical of PIT ratings[3]. As we will see, this will lead to multiple problems.

## 3. The Contradictory Relationship Between the EBA-GL LOM and Supervisory Practices with Internal Ratings.

**The Consistency of EBA-GL LOM with Basel II on the Topic of Ratings.**

The EBA Guidelines on Loan Origination and Monitoring (EBA, 2020; EBA-GL LOM in short) are aligned in many respects with Basel II. The result of ongoing monitoring processes is the potential inclusion of the credit position in a watchlist, which means listing it for future 360-degree review(4); instead, the result of "periodic review" processes is to "review and update any internal ratings/credit scoring" (point 257 of the EBA-GL LOM).

In the EBA-GL LOM, "regular review" or "periodic review" (regular credit reviews of borrowers) is a thorough re-evaluation of borrowers, similar to the initial credit granting process. It aims to "identify any changes in their risk profile, financial position, or creditworthiness compared to the criteria and assessment made at the time of the loan granting, as well as review and update any internal ratings/credit scoring" (point 257 of the EBA-GL LOM). Furthermore, during periodic review, "in addition to monitoring credit and financial metrics, institutions should take into account information related to qualitative factors that may significantly influence the repayment of a loan. Such factors may include information on the quality of management, agreements/disagreements between owners, the structure and flexibility of costs, trends, the size and nature of investments and research and development expenses, as well as the distribution between debt holders and managers (servicers) within the group consolidation" (point 265).

Although the EWI (Early Warning Indicators) within ongoing monitoring require banks to consider a wide range of information(5), the goal of these processes is not to change the rating but rather to place the position on a watchlist. Only after a subsequent 360-degree evaluation can the rating be modified.

In the context of credit granting/periodic review processes, the EBA-GL LOM strongly instructs banks to analyse, both for small and micro enterprises as well as medium-sized and large enterprises, "the debtor's business model and strategy" (points 121 and 144 for the two sectors), "the realism and reasonableness of financial projections" (points 129 and 151), "the feasibility of the business plan and associated financial projections" (points 134 and 161), and the future profitability of the client (points 120 and 152), "under potentially adverse conditions" (points 131 and 156).

Therefore, the EBA-GL LOM requests a comprehensive business analysis in which financial and strategic profiles, both quantitative and qualitative, are integrated. The profitability of the business is the cornerstone of creditworthiness, the time horizon is long-term, and the assessment is made using sensitivity analysis under stress conditions. All this: a) aligns with the long-term orientation of ratings required by Basel II; b) allows for greater allocative efficiency in the banking sector, enabling banks to allocate resources to the most deserving medium-term entrepreneurial initiatives and accompany businesses in investment paths whose returns may not be observable in the short term; c) should be of interest to individual banks for a variety of technical and strategic reasons (a brief summary is in note(6)).

**EBA-GL LOM and Ratings: Contradictions with Basel II.**

In contrast to the clear and consistent framework outlined above, the EBA-GL LOM also contains significant contradictions with Basel II.

The first contradiction emerges in point 274 of the EBA-GL LOM and concerns the purpose of ratings: the internal rating becomes one of the 19 elements to be considered in ongoing monitoring! In fact, within the EWI of ongoing monitoring, banks are asked to consider, among the "indicators of deterioration of credit quality... an actual or expected downgrade of the internal credit rating/credit risk classification for the transaction or client" (point 274, letter q). Here, the internal rating becomes an early warning signal to be used in ongoing monitoring processes, with the goal of placing anomalous positions on the watchlist to later examine them in detail and possibly revise the internal rating. The contradiction is evident. It presupposes an autonomy of the internal ratings and a short-circuit between processes!(7)

A second contradiction arises concerning the scope of information included in the ratings within credit granting/review processes. In EBA, 2020, points 121 and 144 (identical but the former referring to small and micro enterprises and the latter to medium-sized and large enterprises), credit scoring/rating is one of five elements to consider: "In assessing creditworthiness, institutions should: a. analyse the client's financial position and credit risk, as outlined below; b. analyse the organizational structure, business model, and strategy of the client, as outlined below; c. determine and assess the client's credit scoring or internal rating, if applicable, in accordance with credit risk policies and procedures; d. consider all financial commitments of the client, including all credit lines, used and unused, with institutions, as well as credit exposures, repayment behaviour, and other obligations arising from taxes or other public authorities or social security funds; e. assess the structure of the transaction, including structural subordination risks and related terms and conditions, such as restrictive clauses, and, where applicable, third-party personal guarantees and the structure of the real guarantee."

These points suggest that the strategic and financial analysis of a company's business model and the evaluation of financial plans are separated from and additional to the rating. Therefore, the EBA-GL LOM also enshrines in regulation what had already been realized in supervisory practices (as we will see shortly): a) ratings no longer contain all the relevant information that the bank uses to assess debtors/transactions; b) ratings are no longer the final summary of the bank's creditworthiness assessment for credit decisions; c) the rating decision is separated from the loan underwriting decision, creating the split that Basel II had carefully sought to avoid.

**The EBA-GL LOM, Credit Pricing, and Ratings: A Managerial Trap?**

Chapter 6 of the EBA-GL LOM is dedicated to credit pricing. The reason why the EBA guidelines address a topic that might seem to fall within the bank's management autonomy is clear when considering the risk of adverse selection, which we introduced earlier in this article.

Such risk can undermine the stability of financial institutions, so it is of primary concern for supervisors to prevent banks from being exposed to it. Thus, differentiating credit prices according to risk is essential. The problem is: should credit risk be measured using debtors and transactions ratings, or can it be measured using more top-down approaches (such as average loss for product type)?

Key points from the LOM in this regard are as follows (EBA, 2020): "199. ... Institutions should also define their pricing approach based on the type of client and credit quality, and, if applicable, based on the client's risk (in the case of individual pricing determination) ...". "200. Institutions should consider differentiating pricing frameworks based on the type of loan and client. For consumers, micro-enterprises, and small enterprises, pricing should be based more on the portfolio and products, whereas for medium and large enterprises, it should be more closely related to the transaction and the loan." "202. Institutions should consider, and incorporate in loan pricing, all costs ... a. cost of capital (considering both regulatory and economic capital), which should result from capital allocation according to established divisions, such as by geography, business line, and product; ... d. credit risk costs calculated for different homogeneous risk groups, taking into account past loss recognition experience for credit risk, and if applicable, using models for expected loss."

It follows that the EBA's expectations are that for consumers, micro-enterprises, and small enterprises, pricing should be based on historical portfolio and product loss experiences rather than on individual debtor and transaction ratings. This is indeed how the market works, sometime constrained by legal requirements not to discriminate loan prices for different borrowers.

The problem is that when the market price is standardized, banks and financial companies suffer risk-adjusted economic losses when lending to higher-risk clients/transactions, and have an interest in seeking clients/transactions that generate superior risk-adjusted returns due to their lower-than-average risk (thus improving their profitability and stability over time). In other words, they should still try to use adequate rating/scoring systems to guide acceptance/rejection decisions. If they do so, risk-adjusted performance indicators (mentioned in the EBA-GL LOM at point 203: EVA, RORAC, RAROC, RORWA, ROTA) will be able to signal which clients/transactions to accept or reject.

Therefore, if an institution were to adopt the minimalist approach outlined by the EBA-GL LOM for consumer, micro-enterprise, and small business markets and assume that adequate rating/scoring systems are not essential competitive tools, it would expose itself to the risk of adverse selection.

**Contradictions with Basel II in Supervisory Practices.**

A first contradiction arises from the convergence of interests that has led to nearly identifying ratings with statistical-based assignment systems, thus diverging from what is outlined in the previously mentioned Basel II and CRR regulations. In fact, banks have seen advantages in this type of rating in terms of process standardization, cost reduction, and time savings, while supervisory authorities have considered them useful for preserving the integrity and objectivity of evaluations that determine the minimum capital requirements of banks using the Internal Rating Based approaches. Thus, the regulatory validation process of rating systems by the competent supervisory authority has, in fact, required that the rating systems to be validated (and used in management, because of the "use test" previously mentioned) be based on robust statistical models.

A second contradiction with Basel II arises from another convergence of interests between banks and the supervisory authorities in charge of validating the models (initially national supervisors in Europe, before the Single Supervisory Mechanism was operational). Authorities agreed to validate models that violated paragraph 414 of Basel II, discussed earlier, and used a one-year target horizon in the rating assignment phase (not only in the quantification phase). In constructing rating assignment models, the target forecast horizon determines the "observation period" of the dataset in which it is surveyed whether debtors have defaulted or not. The explanatory variables in the model can include all information available at the start of that observation period (usually called "time zero"): internal behavioural data of borrowers' credit lines from the day before, credit register data usually referred to a month and a half earlier, the last approved financial statement (for a small and medium enterprises, it could be referred to a year and a half ago, if time zero falls in the spring months, that is before the new financial statement approval), and qualitative information gathered up to time zero (via the "qualitative questionnaire" or other channels).

The "short" observation period (one year) allows banks to have more historical data to include in the model's estimation dataset, increasing the model's discriminatory power and, therefore, the chances of regulatory validation by national supervisory authorities who, interested in having major banks become IRB banks as soon as possible, chose not to consider paragraph 414 of Basel II.

However, this choice has serious implications on which variables are most predictive and are therefore included in the model (or on which "modules" play a predominant role in the final model, when this is constructed using specialized modules - partial models – for specific types of information). Ultimately, this decision has important consequences on the nature and fungibility of the model in the granting/review process or in ongoing monitoring.

In fact, when the target forecast period is limited to one year, the information that statistical procedures find relevant to optimize the model is essentially limited to internal behavioural data, while financial statement information plays a limited role, and strategic information becomes marginal (Cuneo De Laurentis Salis Salvucci, 2016, and more extensively, Aifirm, 2016)(8). The model thus becomes a substitute for early warning models used for ongoing monitoring; and it seemingly even improves their performance due to the fact that the presence of other variables, in addition to internal behavioural data, limits the excess of false alarms that such models typically suffer from. However, when the same tool is used both for credit granting/review processes and for the ongoing monitoring, open contradictions arise, because of the different objectives of the two processes (reiterated, as mentioned above, by the EBA-GL LOM) and because of the different range of information to be examined in the two processes (according to the same EBA LOM, Basel II, and the CRR).

The problem of the overly narrow scope of relevant information, which ultimately characterizes internal ratings predominantly used by banks, has been further exacerbated by the very restrictive attitude initially taken by supervisory authorities regarding the possibility of using overrides.

A third contradiction with Basel II, which adds to the previous two (concerning assignment methodologies of internal ratings and the range of information contained in them), relates to the nature of ratings as the final summary of the creditworthiness of debtors/transactions, to be used simultaneously in both credit granting/review processes and in regulatory capital adequacy

calculations. Essentially, the legitimate concern to prevent the final assignment of ratings from being influenced by parties with interests that may conflict with ratings' robustness, in some jurisdictions has been pushed to the point of excluding even those who have underwriting powers in credit granting(9).

This has two implications: 1) it pushes banks to establish a separate Credit Risk Management function from the Credit department, and an office often referred to as the rating desk, responsible for the final rating decision; 2) it can create discrepancies between the borrowers' creditworthiness perception of underwriters and the rating assigned.

The informational content of the rating and everything that follows (primarily loan provisions and capital) may, at this point, not fully reflect the assessments that lead to granting loans. This changes the nature of the rating and frustrates one of the main objectives of Basel II: to base capital adequacy on the same internal ratings used by banks to make daily credit risk decisions.

## 4. The Dangerous Relationship Between Rating and Basel 3+ / CRR3

The new capital adequacy regulation (which we refer to here as Basel 3+)(10) confirms many of the Basel II provisions mentioned above (which, as we have noted, have later been ignored by other regulations, guidelines, and supervisory practices), leaving open all the contradictions previously highlighted(11).

Moreover, Basel 3+ seems to move towards a weakening of the role of internal ratings, due to the introduction of the Output floor for the overall result of internal models, set at 72.5% of the RWA calculated using the Standard method, as well as due to various other more specific provisions(12). However, the widespread perception that Basel 3+ reduces the role of internal ratings is incorrect. And it should be in the authorities' interest to stress this point.

First, with respect to internal ratings and from a regulatory standpoint: a) impact assessments (EBA, 2024, Table 4) reveal that portfolios treated with the IRB method result in additional Tier 1 Capital savings for all types of banks, and b) the non-use of internal ratings can have consequences on capital requirements under the Second Pillar of Basel regulations.

Second, the management profile of internal ratings remains crucial: a) they are essential to protect the bank from the risk of adverse selection (overall, Basel 3+ increases capital requirements for banks, and therefore increases the need to align interest rates on loans with the costs to be passed on to debtors), and b) internal ratings for private non-large-corporate debtors are much more productive competitive tools than external ratings for many reasons, the first of which is that the latter are simultaneously available to all credit suppliers and thus do not provide an informational advantage useful for making better decisions than competitors (essentially, they do not eliminate the risk of adverse selection compared to those who internally determine ratings with a more accurate and proprietary information spectrum).

Therefore, the perception held by many that Basel 3+ regulators are no longer pressing for banks to implement appropriate internal rating systems is incorrect. It is also dangerous, as it may lead to a relaxation in the development and use of this essential management tool, that protects from adverse selection risk and safeguards bank's stability.

Now, consider the position on external ratings (issued by recognized rating agencies, the External Credit Assessment Institutions - ECAI) as reflected in the European regulation for implementing Basel 3+: in point 13 of the EU CRR3, 2024, we find: "…institutions should be able to refer to credit assessments by nominated ECAIs to calculate the own funds requirements for a significant part of their corporate exposures… the 'European Supervisory Authorities'… should monitor the use of the transitional arrangement and should consider relevant developments and trends in the ECAI market, impediments to the availability of credit assessments by nominated ECAIs, in particular for corporates, and possible measures to address those impediments. The transitional period should be used to significantly expand the availability of ratings for Union corporates. To that end, rating solutions beyond the currently existing rating ecosystem should be developed to incentivize especially larger Union corporates, to become externally rated. In addition to the positive externalities generated by the rating process, a wider rating coverage will foster, inter alia, the capital markets union. … Member States… should assess whether a request for the recognition of their central bank as an ECAI… and the provision of corporate ratings by the central bank for the purposes of Regulation (EU) No 575/2013 might be desirable in order to increase the coverage of external ratings".

Thus, the emphasis is on the use of external ratings for the corporate market, suggesting that new rating agencies and even the central banks of member states could provide them. In this market, there are also small and medium-sized enterprises, and considering that large corporates are already rated by the major international agencies recognized as ECAIs in Europe, the CRR3 indication seems specifically aimed at increasing the availability of external ratings for SMEs.

This hypothesis seems to be in strong contradiction with the effort made by Basel II to encourage the use of internal ratings and poses further dangers. In fact, the risk/reward ratio of external ratings is largely positive when the object of evaluation is not SMEs, but presents many more disadvantages when the object of evaluation is this type of debtor.

The use of external ratings for SMEs can have advantages in terms of cost-effectiveness in debtor evaluation processes, greater attention to information gleaned from big data, and in terms of "discipline effect" on debtors, which can be of particular interest in some countries(13). However, the use of external ratings for SMEs has several important disadvantages: a) it tends to impoverish the informational base regarding clients' strategic fundamentals (acquired at reasonable costs only by those in continuous and territorial contact with the enterprise), leading conversely to an increased role for internal behavioural data and credit risk bureaus' data in the assignment of increasingly point-in-time ratings; b) it limits the role of credit analysts at banks and discourages the development of business analysis skills; c) it eliminates the informational synergies between credit evaluation activities and the bank's commercial activities, which are at the core of both an effective supply of financial and advising services, and an effective forward-looking evaluation of credit risk.

Ultimately, by promoting external ratings for SMEs, the bank-firm relationship shifts towards a transactional logic (where credit is viewed as a commodity to be produced and sold individually at the lowest possible price) rather than a relational logic (where credit is part of a long-term relationship between the enterprise and the bank). Apart from other macroeconomic implications of this shift,

on the regulatory consistency side, it appears misaligned with the range of information to be considered in granting/reviewing processes foreseen by EBA-GL LOM. As we have noted, these guidelines are very much in line with relationship-oriented banking models, and much less with credit assignment processes of institutions oriented towards transactions, asset-based lending, and instant lending.

It should be noted that discouraging the ability to acquire soft information in relationship lending undermines one of the main reasons for the existence of banks, according to financial intermediation theory, and exposes banks even more to new non-bank competitors.

## 5. Conclusions

Despite some contradictions (preference for TTC ratings but request of a good calibration typical of PIT ratings), for the most part Basel II has a very clear and productive view of internal ratings: these are important tools for the adequate management of credit risk in banks and empower their competitiveness; they represent the final synthesis of the creditworthiness assessment on debtors/transactions, to be used both in credit granting/review processes and in regulatory capital adequacy calculations, to be assigned with a medium term target horizon, leaving banks the choice of assignment methods as long as all relevant information (including qualitative and forward-looking info) is included in the them.

Instead, in the EBA-GL LOM and in supervisory practices ratings have become expressions of only a part, sometimes a minority part, of the relevant information for loan underwriting; they are the result of statistical tools where the role of behavioural information is predominant over qualitative, strategic, and forward-looking information; they are often approved by parties other than those who underwrite loans; they target short-term forecast horizons, getting closer to, if not overlapping with, tools used in credit ongoing monitoring processes.

Basel 3+ confirms the Basel II approach, ignores the contradictions in EBA-GL LOM and the actual evolution of supervisory practices, and generates the false perception that internal ratings are no longer an essential tool for protecting banks against adverse selection risk and safeguarding their stability. In addition, CRR3 emphasizes the use of external ratings, encourages the creation of new rating agencies, and allows central banks of member states to become ECAIs, thus pushing towards a commodity-oriented logic at the expense of the relationship-oriented logic of the bank-enterprise relationship (which the EBA-GL LOM seem to deeply inspire).

"In conclusion, a key question arises: Is the evolution of the role and content of ratings in regulatory and supervisory practices the result of a deliberate, conscious change in approach by well-coordinated authorities (even if never explicitly stated)? Or is it the outcome of various positions that, at different times, dominate the drafting of rules, guidelines, and supervisory practices, without any central authority — even a 'Czar of ratings'— ensuring their coherence?"

## Bibliography

Aifirm, 2016: AIFIRM Position Paper On Validation Of Rating Models' Calibration, a cura di Silvio Cuneo, Giacomo De Laurentis, Fabio Salis, Fiorella Salvucci

BdI, 2006: Banca d'Italia, Circolare 263, 27 dicembre

BIS, 2004: Basel Committee on Banking Supervision, June 2004, International Convergence of Capital Measurement and Capital Standards. A Revised Framework

BIS, 2005: Basel Committee on Banking Supervision, May 2005, Studies on validation of internal rating systems, WP n. 14

BIS, 2017: Basel Committee on Banking Supervision Basel 3+: Finalising post-crisis reforms, December 2017

Cuneo De Laurentis Salis Salvucci, 2016: Cuneo S., De Laurentis G., Salis F., Salvucci F., Validation of rating models calibration, Risk Management Magazine AIFIRM, n. 1

De Laurentis Maino Molteni, 2010: De Laurentis G., Maino R., Molteni L., Internal ratings. Methodologies and case studies, Wiley

EBA, 2016: Final Draft Regulatory Technical Standards on the specification of the assessment methodology for competent authorities regarding compliance of an institution with the requirements to use the IRB Approach in accordance with Articles 144(2), 173(3) and 180(3)(b) of Regulation (EU) No 575/2013

EBA, 2017: Guidelines on Pd estimation, Lgd estimation and the treatment of defaulted exposures, November

EBA, 2020: Guidelines on Loan Origination and Monitoring, May

EBA, 2024: Basel 3+ Monitoring Exercise Results based on data as of 31 December 2023, 7 October

ECB, 2019: European Central Bank, ECB Guide to Internal Models, October

EU CRR, 1013: Regulations Regulation (EU) No 575/2013 of the European Parliament and of The Council, 26 June 2013 on prudential requirements for credit institutions and investment firms and amending Regulation (EU) No 648/2012

EU CRR3, 2024: Regulation (EU) 2024/1623 of the European Parliament and of the Council of 31 May 2024 amending Regulation (EU) No 575/2013 as regards requirements for credit risk, credit valuation adjustment risk, operational risk, market risk and the output floor

## Notes

(1)	Paragraph 444 of Basel II: "Internal ratings and default and loss estimates must play an essential role in the credit approval, risk management, internal capital allocations, and corporate governance functions of banks using the IRB approach." BdI, 2006, p.

71: "The rating system is not just a tool for calculating capital requirements, but must play an important managerial role... Banks may only be authorized to adopt the internal ratings-based approach for calculating capital requirements if the rating system plays a vital role in credit granting, risk management, internal capital allocation, and governance functions of the bank." EBA, 2017: "The concept of the use test was introduced in the IRB Approach to ensure the high quality of risk parameters, assuming that institutions would not use estimates of risk parameters for internal risk management unless they believed these estimates appropriately reflect the actual level of risk."

(2) In the PIT logic, the assignment of ratings and their quantification in terms of PD are focused on shorter time horizons, based on the current conditions of the debtor and the sector/economy. During a recession, ratings typically migrate en masse to worse rating classes; therefore, an increase in risk is reflected in a deterioration of the rating; conversely, when the economic or sectoral cycle improves. In addition, PDs are associated with rating classes based on averages of default rates from shorter periods, closer to the current economic situation; as a result, when comparing the PDs associated with a rating class with the actual default rates observed in subsequent periods for that class, the actual default rates are very close to the expected rates incorporated in the PD estimated for that class. Therefore, each rating class is strongly tied to a particular level of default rates (good calibration).

(3) The main articles of Basel II (BIS, 2004) that explicitly push towards the TTC approach are as follows: a) 414: "Although the time horizon used in PD estimation is one year (as described in paragraph 447), banks are expected to use a longer time horizon in assigning ratings." b) 415: "A borrower rating must represent the bank's assessment of the borrower's ability and willingness to perform contractually despite adverse economic conditions or unexpected events. For example, a bank may base rating assignments on specific, appropriate stress scenarios." c) 447: "PD estimates must be a long-run average of one-year default rates for borrowers in the grade." d) 461: "Banks must use information and techniques that appropriately account for long-run experience when estimating the average PD for each rating grade." e) 463: "Irrespective of whether a bank is using external, internal, or pooled data sources, or a combination of the three for its PD estimation, the length of the underlying historical observation period used must be at least five years for at least one source. If the available observation period spans a longer period for any source, and this data is relevant and material, this longer period must be used." This last point was further emphasized in EBA, 2016, p. 25: "It is desirable that PD estimates are relatively stable over time to avoid excessive cyclicality in own funds requirements. To achieve that, PD estimates should be based on the long-run average of yearly default rates. In addition, as own funds should help institutions survive in times of stress, risk estimates should account for the possible deterioration in economic conditions even in times of prosperity." Despite these repeated calls for ratings characterized by the TTC philosophy, Basel II then defines default by giving great emphasis to the debtor's illiquidity through the 90 days Past Due rule (452: "A default is considered to have occurred with regard to a particular obligor when either or both of the following events have occurred... The obligor is past due more than 90 days on any material credit obligation to the banking group"), which pushes the construction of models towards a PIT philosophy. Furthermore, Basel II requires that the actual default rates observed for each rating class in individual years do not deviate too much from the long-term historical average that identifies the PD (501: "Banks must regularly compare realized default rates with estimated PDs for each grade and demonstrate that the realized default rates are within the expected range for that grade"): this implies the good calibration typical of PIT ratings. This was further specified in WP No. 14 (BIS, 2005), which outlines calibration tests (binomial, chi-square, normal, and the traffic-light approach). For more details on these issues, see Cuneo, De Laurentis, Salis, Salvucci, 2016.

(4) For instance, in point 272 (EBA, 2020), it specifies: "On identifying a triggered EWI event at the level of an individual exposure, portfolio, sub-portfolio or borrower group, institutions should apply more frequent monitoring and, when necessary, consider placing them on a watch list".

(5) Among other things, it mentions: "negative macroeconomic events (including but not limited to economic development, changes in legislation and technological threats to an industry) affecting the future profitability of an industry, a geographical segment, a group of borrowers or an individual corporate borrower, as well as the increased risk of unemployment for groups of individuals…changes in the conditions of access to markets, a worsening in financing conditions or known reductions in financial support provided by third parties to the borrower" (point 274 of EBA-GL LOM).

(6) Firstly, banks also provide medium-to-long-term financing: risks of such financing cannot be adequately covered merely by acquiring real or personal guarantees. Secondly, even formally short-term bank financing, including revocable credit lines, are largely intended to meet the durable financial needs of ongoing businesses, as current assets are actually intended to produce the liquidity necessary to repay loans only if the business ceases operations (they are, in fact, largely permanent current assets). Thirdly, this highlights that even short-term credit lines are really short-term-based only on the assumption that the bank that revokes the credit lines is the first to do so within the broader group of financing banks. In this case, the company will be able to repay the loans to that bank using the available margins in credit lines provided by other banks. However, to achieve this, the bank must make the decision to abandon the customer first: this requires having a superior forward-looking analysis of the client compared to the other banks, meaning that the useful forecasting horizon of ratings is not the formal maturity of the credit lines, but the one capable of anticipating the actions of other banks. Fourthly, the main object of investigation is the client. Certain basic principles of corporate finance apply: "businesses are financed, not individual investments," "there is no direct connection between individual investments and individual corporate loans." In fact, in the case of corporate lending, repayment capacity never directly depends on individual asset items or on the outcome of specific investment operations, but always depends on the overall ability of the business to stay in the market and honour all its financial and non-financial obligations. For these reasons, the recommendation in point 156 of the EBA-GL LOM, which suggests limiting the evaluation horizon to the duration of the loan contract as if the bank had only one loan contract with the debtor, is inappropriate. In the case of revocable credit lines, the evaluation horizon should not be limited to the technical time required to revoke them. Finally, the main subject of creditworthiness assessments should be the client as a whole as, given the high level of competition in banking markets, it is essential to build a portfolio of quality clients who can survive economic cycles and continue to be a source of income for the bank over time. Pursuing a "good operations" policy, instead of focusing on "good clients over time", lays the foundation for the progressive impoverishment of the bank's customer base.

(7) There are no contradictions in the guidelines given by EBA-GL LOM in point 274, letters "p" and "g." In the first case, the reference is to migrations of "aggregates" of internal ratings (letter "p": " negative institution-internal credit grade/risk class migrations in the aggregate credit portfolio or in specific portfolios/segments") and in the second case, it refers to "external

ratings" assigned by agencies or "implied ratings" derived from markets (letter "g": " an actual or expected significant decrease in the main transaction's external credit rating, or in other external market indicators of credit risk for a particular transaction or similar transaction with the same expected life ").

(8)  The situation described regarding short-term rating models, heavily centred on behavioural-based data, is confirmed by the survey contained in the Position Paper of AIFIRM (Italian Association of Risk Managers) on the validation of the calibration of rating models used in the corporate, SMEs, and retail segments. To provide a quantitative indication of the importance of explanatory variables related to behavioural data, credit registers data, financial statements, and others (including qualitative, strategic variables, etc.), banks were asked to express, for each category of data, the ratio between the Auroc of the module dedicated to them and the total Auroc of the final model. The role of data other than behavioural-based and credit registers data is limited, even in the corporate segment.

(9)  For instance, BdI, 2006, Title II, Chapter 1, p. 68.

(10)  BIS, 2017; without an explicit name from regulators, the reform is commonly referred to as the Finalized Basel 3+ post-crisis reforms, Basel 3+, Basel 3+.1, or Basel 3+ Endgame, sometimes called Basel 4. Its implementation in Europe mainly occurred through EU CRR3, 2024.

(11)  For example, paragraph 414 of Basel II (BIS, 2004) we commented on is directly transposed in point 181 of Basel 3+ (BIS, 2017); paragraph 417 in point 185; paragraph 425 in point 193; and paragraph 444 in point 212.

(12)  In EU CRR3, 2024, point 5 states: " The output floor represents one of the key measures of the Basel III reform. It aims to limit the unwarranted variability in the own funds requirements produced by internal models and the excessive reduction in capital that an institution using internal models can derive relative to an institution using the standardised approaches." Among other provisions that limit the role of internal ratings in Basel 3+, we can highlight: IRB methods are prohibited for exposures in capital instruments; for exposures to large companies, banks, and financial sector entities, the regulatory LGD of the FIRBA approach must be used (instead of LGD calculated by internal models, as in the AIRBA approach); the scope of application of internal estimates of CCF/EAD is reduced, new input floors are introduced for PD (from 0.03% to 0.05%), for LGD (differentiated by exposure and collateral), and for CCF/EAD (revolving exposures).

(13)  External ratings, possibly produced by national supervisory authorities (as done by Banque de France) and that consider the level of financial leverage of companies, could stimulate their capitalization; similarly, giving due consideration to ratios between revenues or operating income and assets, or debt, or financial liabilities may encourage a reduction in tax evasion.

# GRC4
## by Augeos

# L'eccellenza nella governance dei rischi. A regola d'arte.

**1 piattaforma, 4 prodotti, +50 moduli.** Dall'**identificazione e mappatura delle informazioni rilevanti** alla **reportistica gestionale, direzionale e normativa**, la nostra piattaforma integra tutto ciò che serve per una governance efficace dei rischi. Grazie a un approccio modulare e flessibile, consente **la valutazione ex-ante ed ex-post dei rischi e dei controlli**, il **monitoraggio attivo tramite KRI e KPI**, e la produzione automatizzata di report gestionali pensati per supportare decisioni strategiche e adempimenti regolamentari. **GRC4 di Augeos** costituisce un unico ecosistema per trasformare la complessità del risk management in un processo fluido, strutturato e sotto controllo. **www.augeos.it**

- Ridotti costi di avvio
- Archivio centralizzato
- Sistema multisocietario
- Visione univoca dei rischi
- Reportistica automatica
- Integrabile con fonti terze
- Scalabile verticalmente e orizzontalmente
- Configurabile e modulare

## Augeos
Services and technology to create value