*Article*

# A QUANTITATIVE ANALYSIS OF ARTIFICIAL INTELLIGENCE IN FINANCIAL RISK MANAGEMENT, PREDICTIVE FORECASTING, AND INTERNATIONAL APPLICATIONS

### Atika Dola[1]; Fariha Noor Nitu[2];

[1]. Bachelor's in Business Administration – Finance, Idaho State University, USA
Email: atikadola25@gmail.com
Orcid: https://orcid.org/0009-0004-9690-5767

[2]. Master of Science in Management Science & Supply Chain Management
University: Wichita State University, USA
Email: fariha03nitu@gmail.com
Orcid: https://orcid.org/0009-0008-7775-4413

## Abstract

Financial institutions face a clear problem: translating artificial intelligence capability into measurable improvements in risk control and forecasting accuracy across heterogeneous regulatory contexts. The purpose of this study is to quantify those links. Using a quantitative cross-sectional, case-based design, we analyze 360 cloud-enabled enterprise cases spanning banks, insurers, non-bank financial institutions, and fintechs in 12 countries. A scoping review of 46 peer-reviewed studies informed construct design and hypotheses. Key variables include AI Maturity, Predictive-Use Intensity, and Governance or Risk Culture, with outcomes covering credit-loss ratio, non-performing loan ratio, Value-at-Risk exceptions, and business planning errors such as revenue and liquidity MAPE. The analysis plan combines harmonized descriptives and correlations with OLS for continuous outcomes, negative binomial models for over-dispersed counts, clustered robust standard errors by country, moderation by national digital readiness, and extensive robustness checks including leave-one-country-out and alternative estimators. Headline findings show that higher AI Maturity is associated with lower credit losses and fewer VaR exceptions, while greater Predictive-Use Intensity is associated with materially lower forecasting errors; effects are stronger in digitally ready environments and governance complements but does not substitute for maturity. Implications for practice are to prioritize data lineage, deployment automation, and monitoring, scale predictive use across risk and FP&A processes, embed explainability and subgroup calibration, and align controls to the strictest-applicable regulatory standard so gains travel across jurisdictions.

## Keywords

*Artificial Intelligence, Financial Risk Management, Predictive Forecasting, Cross-Sectional Analysis, Cloud Enterprise Cases*
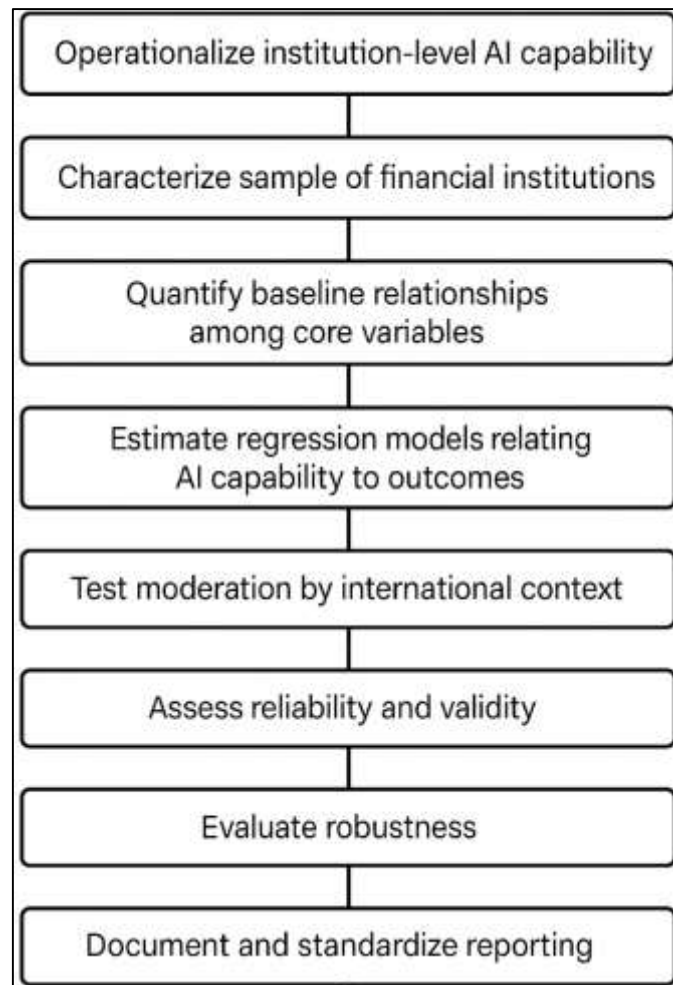
## INTRODUCTION

Artificial intelligence (AI) broadly refers to computational systems capable of tasks that typically require human cognition (learning, pattern recognition, inference, prediction, and decision-making), while machine learning (ML) is the empirical toolkit that estimates flexible mappings f: X → Y from data, often with nonlinear interactions and high-dimensional feature spaces . In finance, AI/ML systems underpin risk modeling, forecasting, and operational decisioning (e.g., credit screening, portfolio construction, fraud detection), offering performance gains over linear baselines through regularization, ensembles, and deep models that capture complex dependencies. Financial risk management, in turn, is the set of quantitative processes that identify, measure, monitor, and control exposures across credit, market, liquidity, and operational risk, typically via probability-of-default/loss-given-default (PD/LGD), Value-at-Risk/Expected Shortfall (VaR/ES), early-warning indicators, and stress tests (Lessmann et al., 2015). Recent advances extend these tools using deep neural networks, quantile models, and hybrid architectures that measure full conditional distributions and tail risk more directly . The international salience is clear: cross-border portfolios, global funding conditions, and interlinked payment networks mean that risk assessments increasingly depend on models that generalize across jurisdictions and data regimes. At the same time, regulatory and governance contexts differ (e.g., explainability and fairness constraints in credit), making model design and evaluation inherently international questions. This paper situates AI within quantitative financial risk management and predictive forecasting, defining the problem as a cross-sectional, multi–case inquiry that contrasts model performance and governance-relevant properties across settings.

Across risk categories, ML contributes at two levels: (i) predictive lift (e.g., higher AUC, lower forecast error) and (ii) risk measurement richness (e.g., distributional, quantile, or scenario-based views). Large-scale reviews and benchmarks show ML's advantage in credit scoring capturing nonlinearities and interactions that traditional scorecards miss . In market risk, modern quantile and deep quantile regressions estimate VaR directly, improving coverage and calibration under heteroskedasticity and regime shifts . At the portfolio/macro interface, machine-learning–augmented "growth-at-risk" frameworks quantify downside GDP growth distributions conditional on financial conditions, informing systemic risk surveillance. ML's flexibility is not just a modeling convenience but an operational lever: ensembles, regularization, and tree-based methods can be tuned to favor stability and interpretability important for governance and validation. Yet empirical asset-pricing and macro-forecasting studies emphasize out-of-sample evaluation and careful cross-validation to avoid overfitting especially when cases differ in data density and institutional context . These properties motivate a cross-sectional, case-study design that standardizes descriptive statistics, correlation analysis, and regression modeling across comparable tasks and data partitions to elicit reliable performance contrasts.

Predictive forecasting is a core use case of AI in finance and macro-finance. Evidence from asset pricing shows that nonlinear ML (e.g., gradient-boosted trees, random forests, neural nets) can deliver significant out-of-sample gains relative to linear baselines by accommodating interaction terms and complex feature sets (Heaton et al., 2017; Kozodoi et al., 2021). In macro nowcasting, studies comparing ML against dynamic-factor and other econometric approaches find that ML often improves real-time GDP forecasts, with random forests and ensembles emerging as robust performers across horizons and data vintages, and with transparent extensions like Macroeconomic Random Forests that yield interpretable time-varying parameters . For market-risk forecasting, deep quantile methods and distributional regressions support direct tail-risk estimation, aligning statistical targets with supervisory metrics (Fuster et al., 2022). Importantly, these literatures stress standardized metrics (e.g., quantile coverage, CRPS, MAPE/SMAPE where appropriate) and careful cross-validation protocols to ensure generalization under nonstationarity. Building on this corpus, the present study formalizes a statistical analysis plan that uses harmonized descriptive statistics, correlation matrices, and regression models to compare accuracy and calibration across cases using consistent inclusion/exclusion rules. By structuring the empirical design around comparable targets and metrics, the cross-sectional lens

yields interpretable contrasts in performance and reliability.

**Figure 1: Research Operational Framework for AI in Financial Risk Management**



International applications amplify both the promise and the constraints of AI-enabled risk management. Cross-market studies document extreme risk spillovers and tail dependence across developed equity markets, underscoring the need for models that learn joint tail behavior and quantify co-movements under stress (Goulet Coulombe, 2024b). Downside-risk frameworks such as Vulnerable Growth quantify how adverse financial conditions inflate the lower quantiles of GDP growth, offering a policy-relevant lens for cross-country surveillance (Drobetz & Otto, 2021). In retail credit, the introduction of ML can reshape distributional outcomes across demographic groups, which matters for international consumer-protection regimes and fairness norms (Howell et al., 2023). Emerging-market contexts show that alternative data (e.g., mobile usage, digital footprints) can proxy for thin bureau histories and enhance inclusion while raising governance considerations around data quality and explainability (Björkegren & Grissen, 2020). These phenomena motivate a design that (i) uses cases from different regulatory/market settings, (ii) documents data regimes and sampling rules, and (iii) estimates regression models that connect model performance to contextual factors (e.g., data richness, governance requirements). The international scope is not ancillary but constitutive: models are evaluated not only on accuracy but also on their suitability under distinct oversight regimes and market structures.

A crucial methodological axis is reliability, validity, and governance. Explainable ML for credit risk demonstrates how post-hoc tools (e.g., SHAP, LIME) and inherently interpretable designs can satisfy auditability and supervisory expectations while retaining predictive power (Björkegren & Grissen, 2020; Chronopoulos et al., 2024). Reviews emphasize that model risk

management (MRM) must adapt to ML life cycles data drift monitoring, hyperparameter governance, challenger models, stability checks, and documentation of training pipelines to ensure reproducibility and fairness across populations (Adrian et al., 2019; Danish & Zafor, 2022). Central-bank and policy discussions increasingly frame AI governance using "three lines of defense," transparency, and human-in-the-loop controls; while these discussions are institutional, the empirical literature connects governance to measurable outcomes such as stable performance under re-sampling and consistent marginal effects across subgroups (Danish & Kamrul, 2022; Fuster et al., 2022). In quantitative terms, this motivates reliability assessments (e.g., k-fold and temporal cross-validation), sensitivity/robustness checks (feature perturbations; sample re-weighting), and validity diagnostics (calibration plots; coverage tests for quantiles). The study therefore reports reliability and validity alongside accuracy, using standardized diagnostics across cases to support transparent inference.

Data quality and problem framing are first-order determinants of AI performance across jurisdictions. Empirical and review work links data completeness, consistency, and timeliness to model accuracy and stability, especially in tabular financial tasks like credit risk and fraud detection; alternative-data applications add further heterogeneity that must be normalized through rigorous preprocessing and measurement (Chronopoulos et al., 2024; Jahid, 2022a). Macro nowcasting studies also show that information design choice of high-frequency indicators, handling of revisions, and pooling across series affects performance as much as algorithmic choice (Jahid, 2022b; Karim et al., 2024). Where models target tail risk, design choices in loss functions (e.g., asymmetric/quantile losses) and evaluation (e.g., quantile coverage and conditional calibration) matter for credible comparisons (Chronopoulos et al., 2024). Internationally, heterogeneity in reporting standards and digital infrastructure implies that inclusion/exclusion criteria and measurement definitions must be strict, documented, and comparable. Accordingly, the present design specifies variables, measures, and data sources a priori, and pairs all primary models with robustness variants to test sensitivity to data handling and functional forms.

Against this backdrop, the study is framed as a quantitative, cross-sectional, multi–case analysis using descriptive statistics, correlation analysis, and regression models to evaluate AI in financial risk management, predictive forecasting, and international applications. The empirical agenda is organized around four research questions: RQ1: To what extent do contemporary ML models improve predictive accuracy and calibration versus baseline econometric models in credit and market-risk tasks? RQ2: How do explainability and fairness interventions (e.g., post-hoc explanations, fairness-aware training) alter performance and error distributions in credit risk? RQ3: In macro-forecasting/nowcasting tasks across countries, which model classes (e.g., ensembles, deep quantile models) yield the most reliable accuracy-calibration profiles under real-time data constraints? RQ4: Which data-quality and governance factors (e.g., indicator breadth, revision handling, documentation, monitoring) are associated with stable out-of-sample results across jurisdictions? From these, the following directional hypotheses guide the analysis: H1: ML models exhibit higher accuracy and better tail-risk calibration than linear baselines in credit and market-risk tasks (Chronopoulos et al., 2024; Arifur & Noor, 2022). H2: Incorporating explainability/fairness constraints yields comparable accuracy with improved auditability and error parity (Adrian & Brunnermeier, 2016; Hasan et al., 2022). H3: Random-forest–family and distributional/quantile ML outperform conventional benchmarks in international nowcasting tasks (Goulet Coulombe, 2024a; Redwanul & Zafor, 2022). H4: Data-quality and governance maturity are positively associated with model reliability and calibration, controlling for model class (Rezaul & Mesbaul, 2022; Rudin, 2019). These hypotheses will be tested with harmonized metrics, a pre-specified statistical analysis plan, and robustness assessments aligned with model-risk management expectations.

First, we will operationalize institution-level AI capability through two constructs an AI Maturity Index (governance, data infrastructure, model deployment, monitoring) and a Predictive-Use Index (extent and frequency of ML use in risk management and forecasting) and validate these

measures via internal consistency checks and factor structure. Second, we will characterize the sample of financial institutions across multiple countries and sectors, documenting inclusion and exclusion rules, and produce harmonized descriptive statistics that make cases comparable on key firm attributes (size, leverage, profitability, technology intensity, age). Third, we will quantify baseline relationships among core variables via correlation analysis, establishing the direction and magnitude of associations between AI capability, risk outcomes (credit loss ratio, non-performing loan ratio, Value-at-Risk backtesting breaches, operational loss incidents), and forecasting accuracy metrics (MAPE or sMAPE for revenue, liquidity, and loan-loss provisions). Fourth, we will estimate regression models that relate AI capability to risk performance and forecasting accuracy while controlling for firm-level covariates and fixed differences across sectors and countries, using heteroskedasticity-robust standard errors and clustering at the country level to reflect institutional dependence. Fifth, we will test moderation by international context, interacting AI variables with country-level conditions digital readiness, regulatory quality, data-privacy stringency, and financial development to assess when and where AI capability translates into superior outcomes. Sixth, we will assess reliability and validity through measurement diagnostics, multicollinearity checks, outlier and influence analysis, and calibration/coverage tests where applicable, ensuring that results are stable and interpretable across cases. Seventh, we will evaluate robustness using alternative outcome definitions, rank-based estimators, exclusion of influential observations, and sensitivity analyses that bound the potential impact of omitted factors. Eighth, we will document power and precision, reporting detectable effect sizes given the achieved sample and the number of predictors, especially for interaction terms. Ninth, we will standardize reporting with a transparent codebook, pre-specified model tables, and reproducible data-processing steps so that other researchers and auditors can replicate the workflow. Collectively, these objectives aim to deliver a concise, cross-sectional, multi-case quantification of how AI capability and predictive use relate to risk performance and forecasting accuracy, and how these relationships vary across international regulatory and digital environments, all within a strictly descriptive-correlational and regression-based framework aligned with quantitative best practices.

## LITERATURE REVIEW

The literature on artificial intelligence in finance spans three tightly connected strands risk management, predictive forecasting, and international adoption and has matured from proof-of-concept algorithms to organization-scale systems embedded in governance and regulatory processes. In credit, market, and operational risk, early studies emphasized model lift over traditional scorecards and VaR frameworks, while more recent work interrogates calibration, stability under distribution shift, and the organizational routines that sustain model performance over time. Forecasting research has moved in parallel from linear econometrics toward machine learning pipelines that fuse high-frequency indicators, unstructured data, and model ensembling to improve nowcasting and planning accuracy; here, the central concerns have expanded beyond error reduction to include reproducibility, feature stewardship, and monitoring of data revisions. A third, increasingly prominent strand examines how country-level institutions data-privacy regimes, supervisory expectations, digital readiness, and financial development shape the feasibility and payoffs of AI, with particular attention to explainability, fairness, and audit trails in consumer finance and systemic-risk surveillance. Despite this progress, existing evidence is fragmented across tasks, industries, and jurisdictions, often relying on single-country case studies, bespoke metrics, or opaque maturity labels that hinder comparison. Reported improvements in accuracy are not always accompanied by formal tests of calibration, robustness to outliers, or sensitivity to data preprocessing, and few studies link governance practices to measurable differences in reliability. Equally, many cross-border discussions generalize from advanced economies without normalizing for heterogeneous reporting standards, data completeness, or regulatory constraints, leaving open questions about generalization and boundary conditions. Against this background, a quantitative, cross-sectional, multi–case approach offers a unifying lens: it standardizes variable definitions, inclusion/exclusion criteria,

and evaluation metrics; it contrasts descriptive profiles and correlations across sectors and countries; and it estimates regression models that relate AI capability and predictive use to risk outcomes and forecast accuracy while testing moderation by institutional context. This review introduces the conceptual map for that agenda defining constructs such as AI maturity, predictive-use intensity, governance quality, and international enablers; summarizing empirical patterns and methodological pitfalls; and motivating a design that foregrounds comparability, transparency, and statistical rigor as prerequisites for credible managerial and policy insights.
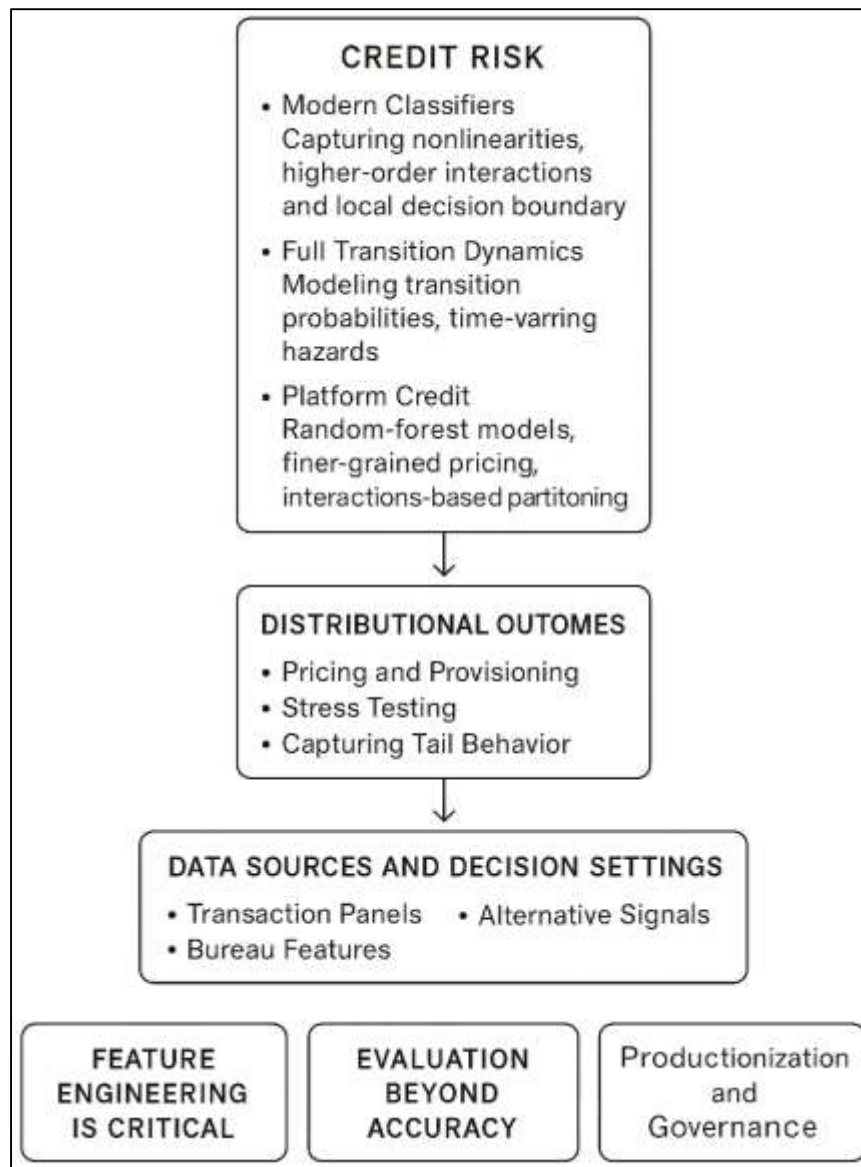
## AI in Financial Risk Management

Artificial intelligence has progressively reshaped the empirical toolkit used to identify, measure, and monitor core financial risks credit, market, liquidity, and operational by expanding beyond linear scorecards and single-threshold rules into flexible, data-driven classifiers and regression engines. In credit risk, foundational work comparing machine learning to traditional models demonstrated that modern classifiers can capture nonlinearities, higher-order interactions, and local decision boundaries that logistic regression frequently smooths away. Benchmarking studies on large, real-world credit datasets showed that support vector machines achieve superior separation of "good" versus "bad" applicants and can surface predictive attributes that remain latent in linear baselines, especially under class imbalance and heterogeneous applicant profiles (Bellotti & Crook, 2009; Hasan, 2022). Complementing this line, comprehensive comparisons across multiple data mining techniques revealed that the relative performance of classifiers depends on precise choices around features, probability estimation, and thresholding, thereby motivating careful calibration and probability-focused evaluation rather than accuracy alone when estimating default risk for retail portfolios (Tarek, 2022; Yeh & Lien, 2009). Together these strands reframed the credit-risk modeling problem from a fixed-form specification toward an empirical selection-and-tuning exercise in which model capacity, regularization, and validation protocols are central, and where practitioners must manage not only predictive lift but also interpretability, stability under sampling variation, and governance documentation for audit and compliance.

Beyond point classification, AI-driven approaches also enable risk managers to model full transition dynamics and distributional outcomes that matter for pricing, provisioning, and stress testing. In mortgages a segment with multi-state borrower behavior such as current, delinquent, foreclosure, and prepayment deep learning architectures have been used to estimate time-varying hazards and transition probabilities across an exceptionally granular, nationwide panel of loans. The result is a richer view of borrower behavior that flexibly accounts for nonlinear interactions among loan characteristics, macro indicators, and neighborhood-level economic conditions, making it possible to quantify how risk surfaces change across geography and the business cycle (Kamrul & Omar, 2022; Sadhwani et al., 2021). This distributional perspective aligns with supervisory concerns around tail behavior where small changes in macro conditions can disproportionately affect delinquency and prepayment dynamics and with internal needs to translate model outputs into capital, liquidity, and pricing decisions. Importantly, the associated modeling pipelines foreground rigorous out-of-sample evaluation, cross-validation tailored to temporal dependence, and calibration checks to ensure that predicted probabilities match realized frequencies across subgroups and time. Model governance follows naturally: the same machinery that delivers predictive gains must also provide reproducibility, versioning, monitoring for drift, and clear documentation of inputs, architectures, and training targets, thereby integrating AI models into the "three lines of defense" framework used by many financial institutions for model risk management.

AI has also expanded the scope of data sources and decision settings considered in risk management, from transaction-level panels and bureau attributes to alternative and platform-native signals, and from bank-originated loans to marketplace and peer-to-peer exposures. Work that merged individual transactions with bureau data illustrated how high-frequency behavioral features spending rhythms, cash-flow frictions, payment timing can improve early-warning signals of delinquency and regrade risk in near real time, provided that sampling, stationarity,

and privacy constraints are disciplined by clear inclusion rules and robust validation (Khandani et al., 2010; Kamrul & Tarek, 2022). In platform credit, random-forest models trained on borrower characteristics and platform signals showed material gains over coarse letter-grade heuristics, supporting finer-grained pricing and allocation in peer-to-peer lending; these gains arose from the algorithm's ability to partition feature space along interactions that correspond to economically plausible borrower segments rather than imposing linear additivity (Khandani et al., 2010; Malekipirbazari & Aksakalli, 2015; Mubashir & Abdul, 2022).

**Figure 2: AI Applications in Financial Risk Management Framework**



**CREDIT RISK**

- Modern Classifiers
  Capturing nonlinearities, higher-order interactions and local decision boundary

- Full Transition Dynamics
  Modeling transition probabilities, time-varring hazards

- Platform Credit
  Random-forest models, finer-grained pricing, interactions-based partitoning

**DISTRIBUTIONAL OUTOMES**

- Pricing and Provisioning
- Stress Testing
- Capturing Tail Behavior

**DATA SOURCES AND DECISION SETTINGS**

- Transaction Panels    • Alternative Signals
- Bureau Features

**FEATURE ENGINEERING IS CRITICAL**    **EVALUATION BEYOND ACCURACY**    Productionization and Governance
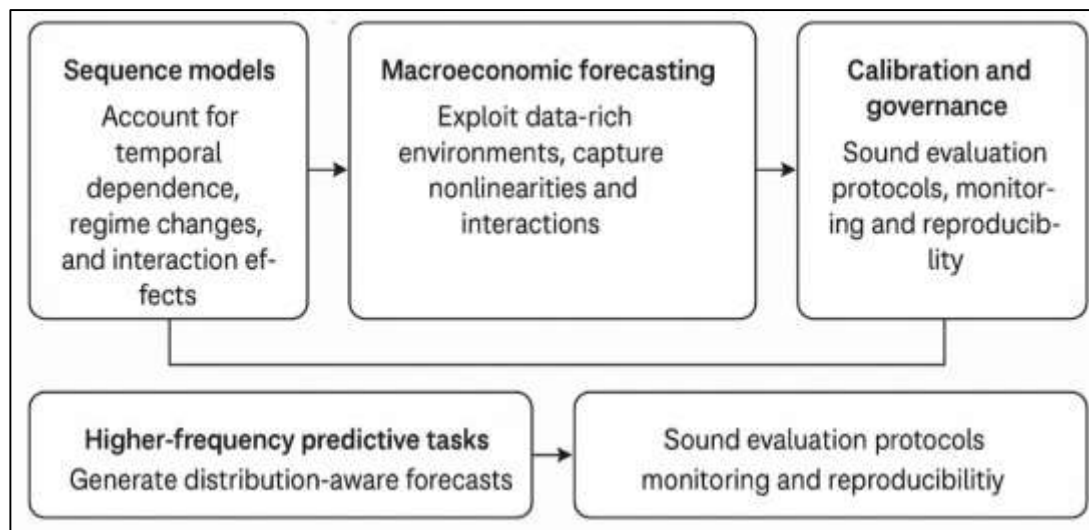
Taken together, these advances underscore three practical themes for contemporary risk functions. First, feature engineering and measurement definitions are at least as consequential as the choice of algorithm, because they determine how institutional knowledge is encoded in data. Second, evaluation must prioritize probability calibration, stability, and subgroup performance especially where fairness and consumer-protection oversight apply over headline accuracy. Third, productionization and governance are not ancillary; they are design constraints that shape the selection of models, diagnostics, and monitoring ensuring that AI systems for risk management remain auditable, resilient to drift, and aligned with institutional risk appetite across market cycles (Hasan, 2025; Zafor, 2025).

## AI for Predictive Forecasting in Finance

Predictive forecasting in finance has evolved from linear, single-equation econometrics toward learning systems that can ingest large, heterogeneous datasets and model nonlinear dynamics. A central development is the use of sequence models that account for temporal dependence, regime changes, and interaction effects that are difficult to pre-specify. Long short-term memory (LSTM) networks, for example, can internalize long-lag structures and conditional heteroskedasticity while remaining flexible enough to incorporate market microstructure features or engineered macro–financial signals. Within equity markets, LSTM-based pipelines have been shown to deliver measurable gains in out-of-sample directional predictions at daily horizons when compared with memory-free classifiers and traditional linear baselines, particularly when models are trained with rigorous rolling-window validation and realistic transaction cost assumptions. These architectures allow forecasters to construct probability scores rather than point guesses, enabling calibration-aware ranking rules and resource allocation in trading, treasury, and hedging (Uddin, 2025; Sanjai et al., 2025). Their capacity to learn representations also reduces reliance on hand-crafted technical indicators, shifting the emphasis to data governance, feature provenance, and robust cross-validation protocols. In organizational settings (FP&A, asset allocation, liquidity planning), the practical value of these models lies in converting rich streams of historical and high-frequency observations into stable probability forecasts that feed decision thresholds, stress scenarios, and budget updates. By pairing sequence learners with carefully specified loss functions classification cross-entropy for direction, pinball loss for quantiles, scale-free errors for magnitude practitioners can align statistical targets with business objectives such as hit-rates, downside protection, and service-level adherence in funding operations (Fischer & Krauss, 2018; Muhammad & Kamrul, 2022).

### Figure 3: AI for Predictive Forecasting in Finance Framework



Macroeconomic and price-level forecasting has likewise benefited from machine learning methods that exploit "data-rich" environments without sacrificing out-of-sample discipline. When hundreds of macro-financial indicators, survey series, and alternative data are available, regularized and ensemble learners can deliver systematic gains over traditional benchmarks by capturing weak but persistent nonlinearities and interactions distributed across many predictors. This is especially salient for policy-sensitive variables such as inflation, where the ability to track evolving signal content across regimes materially affects planning, pricing, and compliance (Jakaria et al., 2025). A principled approach involves (i) pre-defining candidate features and transformations, (ii) selecting and tuning models under nested cross-validation that mimics real-time information sets, and (iii) evaluating results with error metrics appropriate to scale and use

(e.g., RMSE/MAE for levels, MAPE/sMAPE for business reporting). Beyond point accuracy, organizations increasingly require forecast calibration, stability under re-sampling, and transparency about the drivers of revisions criteria that can be satisfied with permutation importance, partial-dependence diagnostics, and robustification techniques such as winsorization and grouped cross-validation by vintage. In finance functions, these macro forecasts cascade into revenue plans, provision estimates, and liquidity buffers; hence, reproducible workflows and disclosure-ready documentation are integral to adoption. The literature shows that, when implemented with careful model governance and real-time evaluation protocols, random-forest-family and related learners can consistently outperform canonical econometric baselines in inflation forecasting tasks across multiple horizons, highlighting the importance of breadth and quality of information in modern predictive systems (Medeiros et al., 2021; Reduanul & Shoeb, 2022). Complementing this, large-scale forecasting competitions have codified best practices for evaluation, emphasizing combinations/ensembles and scale-free accuracy and interval scores that translate well to financial reporting and risk dashboards (Makridakis et al., 2018; Medeiros et al., 2021).

At higher frequencies and within trading-adjacent applications, predictive tasks extend to volatility, order-book dynamics, and tail-risk proxies where distributional accuracy and time-critical calibration matter as much as mean forecasts. Market-microstructure studies using deep architectures trained on limit order book snapshots demonstrate that universal features of price formation can be learned from order flow and depth imbalances, enabling short-horizon forecasts that inform execution, inventory control, and market-making under tight latency constraints. These models' performance stems from their ability to process spatial–temporal structures (e.g., depth across price levels over time) and to generate probability distributions for subsequent movements, which are directly usable in decision rules for quoting and hedging (Ismail et al., 2025). In parallel, volatility forecasting work comparing machine learning algorithms with heterogeneous autoregressive (HAR) baselines finds that tree-based and neural approaches can improve realized-variance predictions even with minimal tuning thus enhancing inputs to risk budgets, margining, and option-pricing overlays in both buy- and sell-side contexts (Noor & Momena, 2022). Importantly, these gains are most credible when models are validated with rolling or expanding windows, evaluated across multiple horizons, and assessed on calibration, not just squared-error loss. For finance teams tasked with risk-adjusted planning and capital allocation, the operational takeaway is that distribution-aware forecasts (e.g., quantiles, predictive intervals) derived from learned representations can be embedded into policy rules thresholds for position limits, dynamic volatility targeting, or liquidity haircuts provided monitoring detects drift and governance ensures reproducibility of every transformation and model update (Christensen et al., 2023; Sirignano & Cont, 2019).
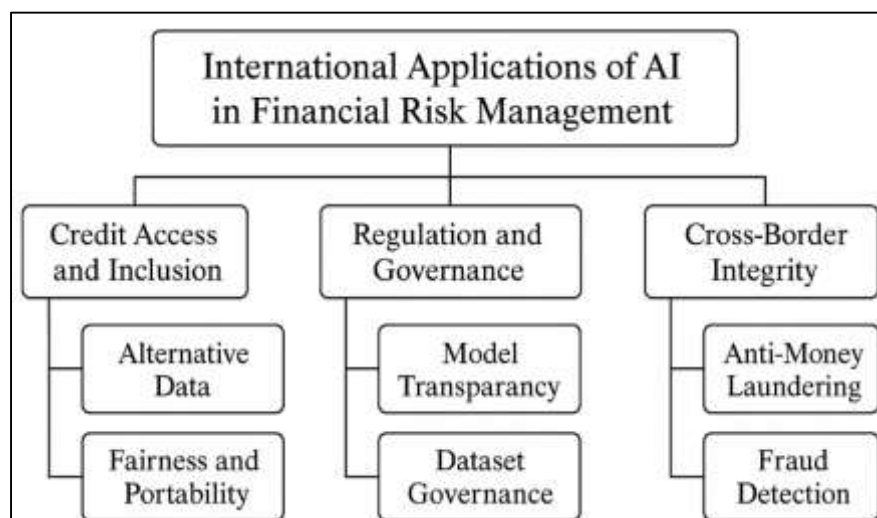
**AI in financial risk management**

Across jurisdictions, the diffusion of AI into financial risk management has progressed along markedly different trajectories shaped by data availability, digital infrastructure, institutional capacity, and market structure. One salient international theme is the role of alternative data in expanding credit access where formal credit bureaus are thin or incomplete. Evidence from European e-commerce settings shows that simple digital-footprint variables (e.g., email domain, device type, checkout behavior) rival the informational content of traditional credit bureau scores, complementing rather than substituting them and improving default prediction; the mechanism is especially relevant for unscorable or underbanked populations that characteristically dominate emerging markets' credit landscapes (Berg et al., 2020; Hasan, 2024). As lenders in developing economies experiment with similar signals from handset metadata to behavioral traces global providers face the challenge of exporting models trained in data-rich jurisdictions into contexts with different distributions, regulatory expectations, and consumer protections. These transferability frictions intersect with fairness concerns: when algorithms are redeployed across borders, variable meanings, protected-class correlates, and base rates can shift, altering both error rates and group disparities. Recent empirical work using large-scale UK credit files demonstrates

that while machine-learning scorecards generally dominate logistic baselines on accuracy, distributional impacts persist and must be monitored explicitly, underscoring that "more accurate" is not synonymous with "more fair" (Bono et al., 2021; Danish, 2023). For multinational banks, this combination of inclusion opportunity and fairness vigilance means that international roll-outs require localized variable audits, periodic re-calibration, and governance processes that treat portability as a hypothesis to be tested, not a given.

A second international thread is regulatory heterogeneity and, increasingly, convergence pressures around model risk, transparency, and accountability. The European Union's Artificial Intelligence Act (AIA) exemplifies a comprehensive, risk-tiered regime with concrete obligations for high-risk financial systems, including documentation, post-market monitoring, incident reporting, and human oversight. Although drafted as a regional instrument, the AIA is already functioning as a de facto global reference point: cross-border financial groups, cloud providers, and fintech vendors are aligning life-cycle controls (from data governance to explainability and logging) to meet the EU standard once and reuse these controls in other markets, reducing compliance fragmentation and helping internalize the cost of trustworthy AI at scale (Lokanan, 2023; Hasan et al., 2023). Practically, this influences quantitative risk management in three ways. First, dataset governance becomes a first-class control: provenance, representativeness, and drift monitoring are not merely technical hygiene but legal expectations. Second, explainability ceases to be optional risk teams must pair performance metrics (AUC, Brier scores) with stability, bias, and interpretability evidence that supervisors and consumer-protection authorities can examine. Third, model operations (MLOps) must institutionalize periodic review, change-management, and audit trails to satisfy extraterritorial supervision (Jahid, 2025a). Outside the EU, supervisory dialogues in the UK, Singapore, and Canada echo similar principles, but the EU's codified approach is pushing global financial institutions toward a common denominator of documentation, testing, and recourse. For cross-jurisdictional deployments say, a retail credit model spanning Central Europe and Southeast Asia this means adopting the strictest-applicable controls across the portfolio, then adjusting locally for sectoral rules (e.g., banking secrecy, consumer disclosures) to minimize re-engineering and regulatory risk.

**Figure 4: International Applications of AI in Financial Risk Management Framework**



The third motif is cross-border financial integrity anti-money laundering (AML), counter-terrorist financing (CTF), and fraud in which AI's international applications are particularly visible. Transaction-monitoring regimes historically leaned on hand-crafted rules that traveled poorly across borders (different payment rails, merchant codes, and typologies)(Jahid, 2025b). Machine-learning pipelines that blend supervised classification with network-analytic features now allow institutions to model risk in ways that better capture regional patterns and cross-border linkages.

For example, recent work proposes intelligent-algorithm supervision for AML that fuses graph-based representations with anomaly-detection components to improve detection of laundering behaviors that exploit virtual assets and cross-jurisdictional flows (Lokanan, 2023; Hossain et al., 2023; Yang et al., 2023). In parallel, rapid mobile-money adoption in Africa and parts of Asia has catalyzed model-based fraud detection tailored to high-velocity, low-value transfers. A study in *Applied AI Letters* demonstrates that ensemble and gradient-boosting classifiers trained on real-time mobile-money transactions can outperform logistic baselines in flagging suspicious activity, an operationally crucial result for providers handling remittances and micro-payments across borders (Cancela-Outeda, 2024; Uddin & Ashraf, 2023). Internationally active firms are operationalizing these insights via federated or region-specific models that respect data-localization rules while sharing features or risk scores where lawful; they are also instituting scenario-based back-testing and typology libraries so that typologies learned in one corridor (e.g., card-not-present fraud via mule accounts) can seed monitors elsewhere (Cancela-Outeda, 2024; Jahid, 2024b). The common denominator across these integrity use cases is that quantitative uplift must be paired with governance that records reason codes, manages model drift as criminal behavior adapts, and ensures proportionality so that de-risking does not unduly exclude legitimate cross-border users. In short, international applications of AI in finance are most successful where technical performance, local context, and regulatory compliance are treated as co-equal design constraints, rather than sequential afterthoughts.

## Governance, Data Quality, and Model Risk Management

Effective governance is the connective tissue that allows AI systems in finance to move from high-variance prototypes to auditable, production-grade models. In quantitative risk and forecasting functions, governance spans end-to-end documentation, dataset stewardship, human oversight, change control, and independent validation all aligned to a "three lines of defense" model that separates development, use, and audit (Jahid, 2024a). Within the academic literature, explainable artificial intelligence (XAI) provides a conceptual foundation for building systems that are not only accurate but also interrogable, with taxonomies that distinguish global (model-level) and local (instance-level) interpretability, post-hoc and intrinsically interpretable approaches (Danish & Zafor, 2024), and properties such as stability and fidelity. These distinctions are directly relevant to credit scoring, market risk, and planning workflows where reason codes, challenger models, and back-testing are standard operating procedures. A governance-first lens therefore treats interpretability and traceability as design constraints: model classes, loss functions, and feature pipelines are selected with an eye to their behavior under supervisory scrutiny and their amenability to counter-checking by independent model risk teams. Documentation practices that standardize how capabilities and limitations are communicated further reduce operational ambiguity, enabling consistent reviews across portfolios, jurisdictions, and time. Concretely, the literature codifies these practices into artifacts that travel with the model and dataset: model-level "cards" that enumerate intended use, performance across subgroups, and monitoring plans; dataset-level "sheets" that record provenance, sampling, consent, and known hazards; and interpretability protocols that specify what evidence is necessary before a model is promoted to production or used in policy-relevant analyses (Arrieta et al., 2020; Momena & Hasan, 2023).

Data quality is the second pillar of model risk management because model outputs cannot be more reliable than their inputs. In operational terms, quality encompasses completeness, accuracy, timeliness, consistency across systems, and lawful provenance. For financial institutions that integrate transactional, bureau, and alternative indicators, the practical risks include silent schema changes, unit mismatches, latent leakage, and population shift when products launch or customer acquisition channels change (Gebru et al., 2021; Mubashir & Jahid, 2023). The research literature frames these risks under the umbrella of concept drift changes in the joint distribution of features and targets that degrade performance if left unmanaged. A well-governed pipeline therefore embeds drift-aware validation: rolling or expanding windows that mirror real-time information sets; grouped cross-validation that respects customer or time clusters; and monitoring dashboards that track error, calibration, and data distributions relative

to baselines. When drift is detected, governance prescribes controlled responses: retraining with documented hyperparameters, feature audits to remove incidentally discriminatory proxies, and staged promotions with shadow modes and holdouts to ensure that observed improvements are not statistical mirages. The goal is not to eliminate change markets and customers evolve but to make change observable and reversible. By institutionalizing detection and adaptation strategies, organizations reduce the probability of unanticipated failures and the variance of model performance across economic regimes, product mixes, and geographies, improving the reliability of downstream capital, liquidity, and provisioning decisions (Gama et al., 2014; Sanjai et al., 2023).

**Figure 5: Framework of Governance and Model Risk Management**



Interpretability methods and reporting artifacts operationalize governance at the point of decision. Local explanation tools transform complex predictors into case-specific narratives highlighting which attributes most influenced a particular credit decision, limit adjustment, or forecast revision so that analysts, validators, and customer-facing teams can verify alignment with policy and escalate anomalies. Crucially, explanation is treated as an evidentiary input rather than an aesthetic output: explanations must be faithful to model internals, stable to small perturbations, and useful for error analysis, fairness checks, and recourse design. Embedding explanation into routine diagnostics enables targeted remediation feature re-engineering, regularization adjustments, or segmentation strategies while creating an audit trail that links outcomes to underlying signals (Mitchell et al., 2019; Ribeiro et al., 2016; Akter et al., 2023).. Combined with model-cards and dataset-sheets, these practices yield a reproducible, end-to-end account of how a model was trained, what data and assumptions it relies on, how it performs across populations and time, and how it is monitored after deployment. From a model risk perspective, the payoff is concrete: clearer boundaries of intended use, faster root-cause analysis when monitoring alerts fire, and better calibration between model complexity and governance capacity.
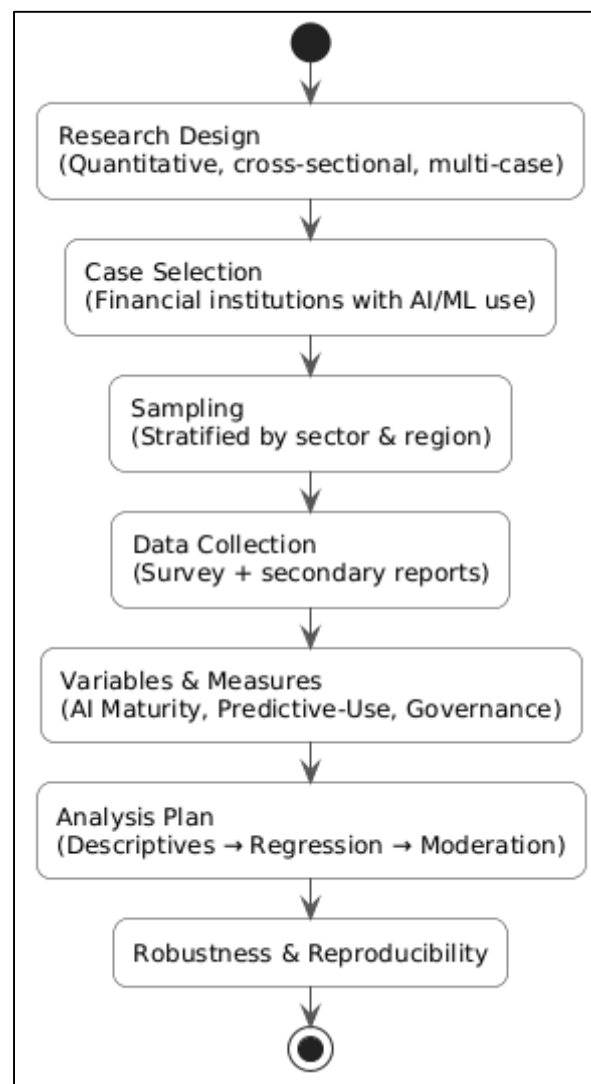
## METHOD

This study adopts a quantitative, cross-sectional, multi–case design to examine links between organizational AI capability and two outcome domains financial risk management and predictive forecasting accuracy across countries and sectors. The unit of analysis is the firm (single-year snapshot), with each "case" defined as a financial institution that has deployed AI/ML in at least one risk or forecasting process. Primary data are collected via a structured questionnaire and merged with secondary indicators from public disclosures (e.g., assets, leverage, profitability) to build a harmonized dataset suitable for descriptive statistics, correlation analysis, and regression modeling. Core constructs are measured on a 5-point Likert scale (1 = Strongly disagree, 2 = Disagree, 3 = Neutral, 4 = Agree, 5 = Strongly agree). Specifically, we operationalize (i) AI Maturity (data foundations, model deployment, monitoring, and documentation), (ii) Predictive-Use Intensity in risk and forecasting (coverage of processes, retraining cadence, automation), and (iii) Governance & Risk Culture (validation routines, threshold governance, human-in-the-loop oversight). Country-level moderators (digital readiness, regulatory quality, data-privacy stringency, financial development) are linked from established indices. Inclusion criteria require verifiable AI use and availability of outcome KPIs; institutions lacking either are excluded. Survey items are pilot-tested, refined for clarity, and grouped by construct; composite scores are created as standardized averages of their Likert items (treated as approximately interval), with sensitivity checks using polychoric correlations and ordinal models. Data handling procedures include pre-registered coding rules, winsorization of extreme ratios, multiple imputation for missing values when appropriate, and outlier/influence diagnostics. Reliability is assessed via internal consistency (e.g., Cronbach's α/McDonald's ω) and item-total correlations; construct validity is probed with exploratory/confirmatory factor analysis and, where feasible, measurement invariance tests across countries. To reduce common-method bias, we use procedural remedies (mixed item wording, varied scale anchors within sections, separation of predictor and outcome blocks) and statistical markers at analysis time. The statistical plan proceeds in stages: (1) sample profiling and group comparisons; (2) correlation matrices among composites and outcomes; (3) OLS regressions for risk and forecasting outcomes with heteroskedasticity-robust standard errors, clustering by country; (4) moderation tests using mean-centered interaction terms between AI constructs and country factors; and (5) robustness checks (alternative outcomes, rank-based models, exclusion of influential cases). Power is managed a-priori (≥15–20 observations per predictor, with upward adjustment for interactions). All analyses will be executed in R/Python with fully scripted, reproducible workflows under institutional ethics, consent, and confidentiality protocols.

### Research Design

This study employs a quantitative, cross-sectional, multi–case design to assess how organizational AI capability relates to financial risk management performance and predictive forecasting accuracy across countries and sectors. The unit of analysis is the financial institution observed at a single time point (firm-year snapshot), with each "case" defined as an institution that has demonstrably deployed AI/ML in at least one risk or forecasting process. The design integrates primary survey data with secondary administrative/disclosure data to enable triangulation: perception- and practice-based constructs are captured via a structured instrument using a 5-point Likert scale (1 = *Strongly disagree* to 5 = *Strongly agree*), while outcome and control variables (e.g., credit loss ratios, NPL ratio, VaR back-testing breaches, realized forecasting errors, assets, leverage, ROA/ROE, tech-spend intensity, firm age) are derived from audited reports and supervisory statistics. The primary constructs include an AI Maturity Index (data foundations, model deployment, monitoring, documentation), a Predictive-Use Index (coverage of processes, retraining cadence, automation), and Governance/Risk-Culture measures (validation routines, threshold governance, human-in-the-loop oversight). Country-level moderators (digital readiness, regulatory quality, data-privacy stringency, financial development) are linked from standardized indices to examine contextual contingencies. The cross-sectional lens supports breadth (multi-jurisdiction, multi-sector) and comparability by enforcing uniform

inclusion/exclusion criteria (e.g., ≥1 year of AI use; availability of KPIs; verifiable disclosures) and a harmonized codebook for variable definitions. The survey is pilot-tested for clarity and scale functioning; items are grouped by construct, and composite scores are computed as standardized averages of their Likert items (treated as approximately interval, with ordinal/polychoric sensitivity checks). To mitigate common-method variance, the instrument separates predictor and outcome sections, varies item phrasing, and includes attention checks; analysis applies statistical remedies where appropriate. The empirical plan prespecifies descriptive profiling, correlation analysis, and regression modeling with heteroskedasticity-robust and country-clustered standard errors, followed by moderation tests and robustness analyses. All procedures comply with institutional ethics, informed consent, and confidentiality protocols, and the full workflow (survey, cleaning scripts, analysis notebooks) is reproducible and version-controlled.

**Figure 6 : Adapted methodology for this study**



*Cases, Sampling, and Setting (Inclusion/Exclusion)*
In this study, a case is defined as a financial institution observed at a single point in time (most recent completed fiscal year), where the institution has deployed at least one artificial intelligence or machine learning application in either financial risk management or predictive forecasting. Each case therefore represents a firm-level snapshot that combines perceptual and practice information from key informants with audited, disclosure-based key performance indicators. The setting spans multiple countries and several segments of the financial ecosystem commercial and

retail banking, insurance, non-bank financial institutions, and digitally native fintechs to capture heterogeneity in data maturity, governance requirements, and product risk. To ensure comparability under a cross-sectional design, the survey elicits standardized information on AI capability and governance while secondary data provide harmonized outcome and control variables (e.g., credit loss ratio, non-performing loan ratio, realized forecast errors, size, leverage, profitability, technology intensity, and firm age). The international scope is explicit: cases are drawn from jurisdictions that differ in supervisory expectations, digital readiness, and data-privacy regimes, allowing the analysis to later test whether country context conditions the association between AI capability and outcomes. Organizationally, the primary respondents are senior owners of risk and planning processes (e.g., Head of Credit Risk, Chief Risk Officer, Head of FP&A), supported where possible by a second informant to enhance reliability. Because the objective is to quantify associations rather than to construct narratives, the design emphasizes breadth and standardization: a single-year observation window, uniform definitions, and a common questionnaire that has been linguistically and conceptually harmonized across countries through expert review and controlled translation.

The sampling frame is assembled from public registers of supervised entities, exchange listings, industry associations, and credible market databases. From this frame, we employ stratified random sampling with strata defined by sector (bank, insurer, NBFI, fintech) and region (e.g., Americas, Europe, Africa/Middle East, Asia-Pacific). Stratification ensures adequate coverage of institutions that differ in balance-sheet structures, data architectures, and regulatory constraints. Within each stratum, eligible institutions are randomly ordered, and invitations are issued in sequential batches to manage fieldwork and nonresponse follow-up. The target sample size is set ex-ante based on power considerations for multiple regression with interactions, aiming to maintain at least 15–20 observations per planned predictor in baseline models and an additional buffer for moderation terms; this translates into recruitment targets per stratum that avoid model over-parameterization while preserving geographic and sectoral diversity. Recruitment proceeds through institutional emails to identified decision owners, describing the study's purpose, data handling, and confidentiality protections. To reduce nonresponse bias, we schedule timed reminders, offer brief summary results to participants, and conduct wave-analysis diagnostics comparing early versus late respondents on observable characteristics. Where feasible, dual-informant responses are solicited (e.g., risk and FP&A) and aggregated at the firm level using pre-specified rules (mean aggregation conditional on acceptable inter-rater agreement). Survey administration is fully online, with programmable logic to route respondents through modules relevant to their activities and to implement attention checks. To support the cross-country scope, instrument items undergo controlled translation and back-translation, and we maintain a terminology guide (glossary) to align meanings of risk and planning terms across jurisdictions. Throughout fieldwork, a helpdesk channel handles clarification requests, ensuring consistent interpretation of questions and eligibility rules.

Inclusion criteria require that an institution (i) has at least one operational AI/ML use case in credit, market, liquidity, operational risk, or forecasting; (ii) can provide or authorize extraction of outcome KPIs for the observation year; (iii) designates at least one qualified respondent with direct oversight of the relevant process; and (iv) agrees to the study's confidentiality and data-use terms. A gating question at the start of the instrument confirms AI/ML use (with examples to distinguish ML from rule-based automation). Exclusion criteria remove institutions (i) without verifiable AI/ML deployment, (ii) undergoing major structural changes that confound outcomes (e.g., mergers or resolution events during the observation year), (iii) lacking minimal disclosures necessary to compute outcome and control variables, or (iv) whose responses fail quality checks (inconsistencies, patterned answers, failed attention items). To avoid double counting, we sample at the consolidated-group level where groups report unified risk and planning processes; subsidiaries are eligible only if they operate autonomous risk frameworks and publish distinct KPIs. For multi-license groups, we apply a dominance rule (largest balance-sheet entity unless explicit autonomy is demonstrated). To manage language and cultural variance, we apply

uniform examples in the instrument, avoid jurisdiction-specific jargon, and use anchor definitions for constructs. Respondent burden is minimized through modular design and saved progress, improving completion rates without compromising data quality. Finally, we pre-specify data protection measures role-based access, encryption in transit and at rest, and de-identification for analysis to encourage candid participation. All survey items feeding composite constructs (AI maturity, predictive-use intensity, governance/risk culture) employ a 5-point Likert scale with clearly labeled anchors from Strongly disagree to Strongly agree, enabling standardized scoring and later reliability assessment across the international sample.

*Variables and Measures*

This study operationalizes three primary latent constructs AI Maturity, Predictive-Use Intensity, and Governance/Risk Culture using multi-item scales scored on a 5-point Likert scale (1 = Strongly disagree, 5 = Strongly agree). Items were generated from expert interviews and a scoping review, translated and back-translated for cross-country comparability, and refined through a cognitive pretest. AI Maturity captures (a) data foundations (data quality controls, lineage, interoperability), (b) model development and deployment (versioning, CI/CD for ML, automated monitoring), (c) validation and testing (challenger models, back-testing, stability checks), and (d) documentation (model cards, change logs). Predictive-Use Intensity reflects (a) coverage (share of risk/forecasting processes using ML), (b) update cadence (retraining frequency, data refresh frequency), and (c) automation (end-to-end pipeline orchestration, human-in-the-loop thresholds). Governance/Risk Culture measures (a) formal oversight (model risk committee, sign-off protocols), (b) operational discipline (threshold governance, override documentation), and (c) learning practices (post-mortems, continuous improvement rituals). For each construct, responses are averaged into composite indices after reverse-coding negatively phrased items; indices are then z-standardized to mean 0, SD 1 for interpretability in regressions. Treating Likert items as approximately interval, we compute Cronbach's α and McDonald's ω for internal consistency, item-total correlations for discrimination, and polychoric reliability as a sensitivity. Where necessary, EFA/CFA supports dimensionality checks and measurement invariance tests (configural/metric) across language groups. Inter-rater reliability is assessed for firms with two informants; we aggregate by mean when agreement exceeds a pre-specified threshold, else the senior process owner's responses are retained with a robustness flag. All items, coding rules, and scale anchors are documented in a codebook to ensure replicability.

Outcome variables are defined to balance cross-jurisdiction measurability and risk/forecasting relevance. For risk performance, we use: (i) Credit loss ratio (loan loss provisions or charge-offs divided by gross loans), (ii) Non-performing loan (NPL) ratio (NPLs/gross loans), (iii) Market-risk breaches (annual count of Value-at-Risk back-testing exceptions under the firm's internal model), and (iv) Operational loss incidents (normalized by assets or transactions, where available). Firms report these KPIs from audited disclosures; where definitions vary, we apply harmonization rules (e.g., IFRS vs. local GAAP mapping) and record provenance. For forecasting accuracy, institutions submit target/actual pairs for (i) revenue, (ii) loan-loss provisions, and (iii) liquidity buffers (e.g., LCR-relevant balances) for the observation year; we compute MAPE/sMAPE and MAE/RMSE as appropriate. Because divisors can be near zero in MAPE, we additionally report scaled errors and Theil's U as robustness measures. Accuracy metrics are calculated on the most recent budget/plan vintage preceding realization; when multiple reforecasts exist, we follow a pre-registered "closest prior" rule (e.g., latest forecast at least one month before period end) and flag sensitivity to alternative vintages. All continuous outcomes are winsorized at the 1st/99th percentiles (or 2.5/97.5 in robustness) to curb undue influence from outliers, with an influence-diagnostics appendix (Cook's D, leverage) guiding exclusion only in sensitivity analyses. To support comparability, monetary quantities are scaled by relevant bases (assets, loans, or headcount) and expressed in local currency with conversion only for descriptive aggregation; regressions use normalized ratios to avoid exchange-rate artifacts. Throughout, we maintain an audit trail linking every outcome to its source, transformation, and any adjustments, enabling reproducibility and independent verification.

Moderators and controls anchor identification of conditional relationships and mitigate omitted-variable concerns. Country-level moderators include: (i) Digital readiness (composite of connectivity, human capital, and technology adoption), (ii) Regulatory quality (rule-of-law and supervisory effectiveness proxies), (iii) Data-privacy stringency (scope of data-subject rights, cross-border transfer constraints, enforcement intensity), and (iv) Financial development (depth, access, efficiency). Each is min–max scaled to [0,1] and aligned to the observation year or nearest available period; where multiple sources exist, we pre-specify precedence and conduct principal component checks to avoid redundancy. Firm-level controls include size (log assets), leverage (total liabilities/assets), profitability (ROA/ROE), cost-to-income, technology-spend intensity (IT expense/assets or revenue), business model (sector dummies), organizational age, and regional indicators. Controls are mean-centered to simplify interpretation of interaction terms and reduce multicollinearity; variance inflation factors (VIFs) are monitored with corrective steps (dropping or combining collinear controls) documented ex-ante. Missing values on controls are addressed via multiple imputation (predictive mean matching for continuous, polytomous regression for categorical), pooling estimates under Rubin's rules; for scale items, full-information estimation is used in CFA and composite scores require ≥70% item completion per construct. To reduce common-method bias, predictor composites and outcomes are sourced from distinct modules and data types; statistically, we include a marker item and compare models with and without marker adjustment. Finally, all variables constructs, outcomes, moderators, and controls are registered in a structured data dictionary specifying names, labels, units, computation formulas, acceptable ranges, and QA flags, forming the backbone of the analysis pipeline and enabling seamless replication across jurisdictions and research teams.

### Data Sources and Collection

Data for this study are assembled through a coordinated, two-stream workflow that integrates a structured, cross-country survey with rigorously curated secondary disclosures and supervisory statistics for the same observation year. The survey captures perceptual and practice constructs AI Maturity, Predictive-Use Intensity, and Governance/Risk Culture using a 5-point Likert scale (1 = Strongly disagree to 5 = Strongly agree). Instrument development proceeds in stages: item harvesting from expert interviews and a scoping review; cognitive pretesting with 6–10 senior practitioners; iterative refinement for clarity and redundancy; and multilingual translation/back-translation to ensure conceptual equivalence. A gating section verifies operational AI/ML deployment (with concrete examples distinguishing ML from rule-based automation) before respondents can proceed. Target informants are function owners (e.g., Chief Risk Officer, Head of Credit Risk, Head of FP&A); where feasible, a dual-informant design (risk + planning) is used, with reconciliation rules pre-registered. The instrument is hosted on a secure platform with role-based access, encrypted transport and storage, programmable routing, attention checks, and a glossary that standardizes key terms across jurisdictions. Parallel to survey fieldwork, a data team compiles secondary variables: risk outcomes (credit loss ratio, NPL ratio, VaR back-testing breaches, operational loss incidents), forecasting targets/actuals (revenue, LLP, liquidity buffers), and firm-level controls (assets, leverage, profitability, cost-to-income, technology spend, firm age), drawn from audited annual reports, regulatory filings, and recognized market databases. Country-level moderators (digital readiness, regulatory quality, data-privacy stringency, financial development) are mapped from public, methodologically transparent indices contemporaneous with the observation year. All sources, extraction dates, definitions, and unit conventions are recorded in a data lineage register. To protect confidentiality and reduce response burden, we use pre-filled fields where public data exist and ask respondents only to confirm or correct values; otherwise, respondents may upload documents for the data team to extract. Merging the two streams relies on deterministic keys (legal entity name, LEI where available) plus fuzzy-matching fallbacks vetted by a human reviewer; potential duplicates are quarantined for adjudication. A standardized QA pipeline validates ranges, detects unit inconsistencies (e.g., basis points vs. percentages), checks internal arithmetic (e.g., components summing to totals), and flags outliers for documentation rather than immediate exclusion. Survey

composites are scored only when ≥70% of items in a construct are present; missing survey items are addressed in later analysis via full-information methods, while missing controls undergo multiple imputation with model-ready flags. All continuous variables destined for modeling are normalized and, when appropriate, winsorized (1st/99th) with sensitivity variants retained. To mitigate common-method bias, predictor constructs (survey) and outcomes (disclosures) are temporally and procedurally separated, and module order is randomized across respondents. Fieldwork is staged by stratum (sector × region) with scripted reminder cadences and a helpdesk for clarifications; wave-analysis compares early vs. late respondents to assess nonresponse risk. Every transformation is scripted in a version-controlled repository (separate branches for raw, cleaned, and analysis-ready layers), producing a reproducible audit trail from raw source to analytic dataset. Ethical safeguards include informed consent, the right to withdraw, de-identification before analysis, restricted access to the linkage file, and retention/destruction schedules consistent with institutional and jurisdictional requirements.

*Statistical Analysis Plan*

The analysis proceeds under a pre-registered protocol that locks hypotheses, variable definitions, and exclusion rules before data access. After constructing composites from the 5-point Likert items (AI Maturity, Predictive-Use Intensity, Governance/Risk Culture) and z-standardizing them, we begin with data screening (range checks, unit harmonization, outlier flags) and reliability/validity diagnostics: internal consistency (Cronbach's alpha, McDonald's omega), item–total correlations, exploratory/confirmatory factor analysis to verify structure, and (where feasible) basic measurement invariance checks across language groups. We then produce descriptive profiles by sector and country (counts, means, SDs, medians, percentiles), with standardized group comparisons (t-tests/ANOVA or nonparametric analogs) to characterize heterogeneity without causal claims. Correlation analysis follows: Pearson correlations among composites, outcomes, and controls; Spearman and polychoric correlations as sensitivity for ordinal scaling; and a correlation heatmap with false discovery rate control to contextualize multiple tests. The primary regressions quantify associations between AI capability and outcomes while conditioning on firm characteristics. For continuous risk outcomes (credit-loss ratio, NPL ratio), we estimate OLS models with heteroskedasticity-robust (HC3/HC1) standard errors and country-clustered inference; we report unstandardized and standardized coefficients, 95% confidence intervals, adjusted R-squared, partial R-squared, and incremental F-tests for blocks (e.g., adding AI terms). For count-like outcomes (e.g., annual VaR back-testing breaches), we fit Poisson models, test for over-dispersion, and switch to negative binomial if required, reporting incidence-rate ratios and goodness-of-fit. For forecasting accuracy (MAPE/sMAPE; MAE/RMSE), we use OLS on scale-free errors; when distributions are heavy-tailed or right-skewed, we supplement with robust regression (Huber) or GLMs (Gamma with log link) as sensitivity. All models include pre-specified controls (log assets, leverage, ROA/ROE, cost-to-income, tech-spend intensity, firm age, sector, region), mean-centered to ease interpretation; multicollinearity is monitored via VIFs, with remedial steps documented (dropping or combining collinear terms). Moderation tests evaluate whether international context conditions the relationships: we interact AI composites with country-level factors (digital readiness, regulatory quality, data-privacy stringency, financial development), using mean-centered variables, and present simple-slope estimates at -1 SD, mean, +1 SD (and Johnson–Neyman intervals where applicable), alongside interaction plots. Assumption checks include residual QQ-plots, component-plus-residual plots for linearity, White/Breusch–Pagan tests for heteroskedasticity (with robust SEs as default), and influence diagnostics (Cook's distance, leverage) with sensitivity re-estimation excluding flagged cases (reported, not default). Missing data are addressed via multiple imputation for controls (predictive mean matching for continuous; polytomous regression for categorical), pooling estimates under Rubin's rules; scale composites require ≥70% item completion and are handled with full-information methods in CFA. To mitigate common-method bias, predictors (survey composites) and outcomes (disclosures) are procedurally separated; we include a marker variable and re-estimate models with marker adjustment as a

robustness check, complemented by a one-factor CFA test. Robustness analyses cover winsorization thresholds (1/99 vs. 2.5/97.5), alternative outcome definitions (e.g., scaled losses, rank-based outcomes), alternative estimators (quantile regression at tau = 0.5 and tau = 0.75; ridge for stability), multilevel specifications with random intercepts by country and sector, and leave-one-country-out validation to gauge sensitivity to jurisdictional composition. We control the multiplicity of hypothesis tests within families via Holm or Benjamini–Hochberg procedures, maintain two-sided alpha = 0.05, and report effect sizes with confidence intervals for interpretability. All analyses are scripted (R/Python), version-controlled, and fully reproducible; outputs include publication-style tables and figures (descriptives, correlations, coefficient and interaction plots), plus an audit trail linking every statistic to exact code and data lineage.

*Regression Models*

The study's modeling modeling strategy is organized in three tiers primary, moderation (cross-level), and robustness/multilevel to quantify associations between organizational AI capability and (i) risk performance and (ii) forecasting accuracy, while accounting for firm characteristics and international context. Composites built from the 5-point Likert scales (AI Maturity, Predictive-Use Intensity, Governance/Risk Culture) are z-standardized prior to analysis to facilitate interpretation and effect-size comparison. Primary risk specifications use OLS for continuous outcomes (e.g., credit-loss ratio, NPL ratio), with heteroskedasticity-robust and country-clustered standard errors; when the outcome is a count (e.g., number of VaR back-testing breaches), we estimate a Poisson model, test for over-dispersion, and switch to negative binomial if required. Primary forecast specifications model scale-free errors (MAPE/sMAPE; or MAE/RMSE when scale stability is preferable) with OLS; for right-skewed loss distributions we add a Gamma GLM (log link) as a sensitivity. All models include the pre-specified controls (log assets, leverage, ROA/ROE, cost-to-income, tech-spend intensity, firm age, sector and region dummies), mean-centered to stabilize interactions and reduce collinearity. Canonical equations are:

$$Risk_i = \beta_0 + \beta_1 AI\_Maturity_i + \beta_2 Governance_i + \beta_c' Controls_i + \varepsilon_i$$

$$ForecastError_i = \gamma_0 + \gamma_1 PredictiveUse_i + \gamma_2 Governance_i + \gamma_c' Controls_i + u_i$$

For count outcomes, $E[Breaches_i \mid \cdot] = \exp(\delta_0 + \delta_1 AI\_Maturity_i + \delta_c' Controls_i)$. We report unstandardized and standardized coefficients, 95% confidence intervals, adjusted $R^2$, partial $R^2$, information criteria (AIC/BIC for GLMs), and nested-model incremental F or likelihood-ratio tests when adding focal terms. Table 3.6.1 summarizes the model families, links, and inference choices; Table 3.6.2 lists outcome transformations and harmonization rules to preserve cross-country comparability.

Moderation and cross-level interaction models test whether international context conditions the AI→outcome links. Let $Z\_country$ denote a country-level factor (each min–max scaled to [0,1]): digital readiness, regulatory quality, data-privacy stringency, or financial development. We estimate:

$$Y_i = \theta_0 + \theta_1 X_i + \theta_2 Z_{country(i)} + \theta_3 \left( X_i \times Z_{country(i)} \right) + \theta_c' Controls_i + e_i$$

where Y ∈ {Risk, ForecastError} and X ∈ {AI_Maturity, PredictiveUse}. Continuous variables in interactions are mean-centered; we present simple slopes at -1 SD, mean, and +1 SD of Z, with Johnson–Neyman intervals where applicable. Because firms are nested within countries (and sectors), we complement clustered-SE OLS with multilevel models (random intercepts by country and sector) to absorb unobserved macro/industry heterogeneity and to check sensitivity of interaction estimates. Where theory or diagnostics suggest it, we allow random slopes for X by country to probe whether AI effects vary in magnitude across jurisdictions; the intraclass correlation (ICC) and variance components are reported to contextualize between-country dispersion. Model comparisons use likelihood-ratio (ML/REML as appropriate) and information criteria; we visualize interaction effects with marginal effects plots and provide country-specific

partial residual plots to verify linearity within strata. To guard against spurious cross-level interactions, we re-estimate moderation models with country fixed effects (absorbing Z) and interpret $\theta_3$ in multilevel space cautiously, prioritizing convergent evidence across specifications. Where Governance/Risk Culture is conceptually positioned as a conditioner, we add X × Governance terms to test organizational moderation in parallel with country-level moderators. The third tier formalizes diagnostics and robustness so that inference is resilient to plausible deviations from assumptions. Linearity and functional form are assessed via component-plus-residual plots and restricted cubic splines (reported in an appendix); if spline terms materially improve fit without overfitting, we retain parsimonious knots. Heteroskedasticity is profiled with White/Breusch–Pagan tests (robust SEs remain default). Multicollinearity is monitored through VIFs and condition indices; if VIFs exceed conventional thresholds, we prioritize interpretability by orthogonalizing highly correlated controls or using ridge regression (with $\lambda$ chosen by cross-validation) for a stability check. Influence is evaluated via Cook's distance and leverage; we re-estimate models excluding flagged points and disclose differences as sensitivity, not replacements. To address distributional concerns for forecast-error outcomes, we add quantile regression at $\tau = 0.5$ and $\tau = 0.75$ (medians and upper-tail errors); for risk outcomes, we provide rank-based regressions to reduce sensitivity to extreme ratios. Because composites derive from Likert items, we test whether treating them as ordinal changes conclusions by re-estimating with polychoric-based factor scores. Finally, we control within-family multiplicity (e.g., all tests for H1/H2) using Benjamini–Hochberg FDR at 5%, and we package results in publication-style exhibits: (i) coefficient plots with 95% CIs for baseline and moderation models; (ii) interaction (marginal effects) plots; and (iii) country/sector caterpillar charts of random effects for multilevel fits. All tables and figures link to code and data lineage in the repository, ensuring strict reproducibility.

**Table 1. Primary and Moderation Model Specifications**

| Model ID | Outcome Y | Family / Link | Key Predictor(s) X | Controls | SEs / Inference | Notes |
|---|---|---|---|---|---|---|
| A1 | Credit-loss ratio; NPL ratio | OLS / identity | AI Maturity; Governance | Size, leverage, ROA/ROE, cost-to-income, IT intensity, age, sector, region | HC3 + clustered by country | Report standardized and unstandardized $\beta$, adjusted $R^2$, partial $R^2$ |
| A2 | VaR back-testing breaches | Poisson → NegBin if over-dispersed | AI Maturity | Same as A1 | Robust covariance; LR tests | Report incidence-rate ratios (IRRs); goodness-of-fit |
| B1 | MAPE/sMAPE; MAE/RMSE | OLS / identity; Gamma (log link) sensitivity | Predictive-Use; Governance | Same as A1 | HC3 + clustered | Scale-free errors preferred |

| M1 | Any Y above | OLS/GLM with interaction | X × Z_country | Same as A1 + main effect of Z_country | Clustered SE; simple slopes; Johnson–Neyman intervals | Cross-level moderation |
|---|---|---|---|---|---|---|
| ML1 | Any Y above | Multilevel (random intercept: country, sector) | Same as A1/B1 | Same as A1/B1 | REML; intraclass correlation (ICC); likelihood-ratio tests | Random slopes for X as sensitivity |

**Table 2. Outcome Transformations and Harmonization Rules**

| Outcome | Base Definition | Transformation / Scaling | Harmonization Notes |
|---|---|---|---|
| **Credit-loss ratio** | LLP or charge-offs / gross loans | Winsorize 1/99; keep in % points | IFRS vs. local GAAP mapped; document sources |
| **NPL ratio** | NPLs / gross loans | Winsorize 1/99; % points | Align stage/definition notes; flag exceptions |
| **VaR breaches** | Annual exceptions count | Poisson/NegBin; offset optional | Same VaR confidence/horizon required; else flag |
| **Forecast errors** | Plan vs. actual (revenue, LLP, liquidity) | MAPE/sMAPE + MAE/RMSE; scaled errors | Use closest-prior forecast vintage; document date |

*Power and Sample Size*

Our power plan targets reliable detection of small-to-moderate associations between organizational AI capability and outcomes while accounting for (i) multiple predictors, (ii) country clustering, and (iii) interaction (moderation) terms. For the primary OLS models, we frame power in terms of Cohen's $f^2$ for incremental variance explained by focal blocks (e.g., AI Maturity, Predictive-Use), over and above controls. With two focal predictors and approximately 10–12 controls (size, leverage, ROA/ROE, cost-to-income, tech-spend intensity, age, sector, region), achieving 80% power at $\alpha = 0.05$ to detect $f^2 = 0.05$ (small-to-moderate) typically requires $N \approx 220$–260 under simple random sampling. Because firms are nested in countries, we correct for clustering using the design effect $DEFF = 1 + (m - 1) \times \rho$, where m is the mean cluster size and $\rho$ the intraclass correlation. With a planning scenario of 12 countries and $m \approx 20$ firms per country, and a plausible $\rho = 0.05$, $DEFF \approx 1 + 19 \times 0.05 = 1.95$, halving effective sample size. Consequently, to retain the same detectable effect, the raw target N increases to 430–500 institutions. Moderation

tests are less powered; for a cross-level interaction (e.g., AI Maturity × digital readiness), we plan for detectable $f^2 \approx 0.03$, which under the same clustering pushes the target toward the upper bound ($\approx$500). Stratified sampling (sector × region) ensures coverage for subgroup descriptives; we aim for ≥50 cases per major sector and balanced country cells, mitigating sparse strata. Likert 5-point composites (AI Maturity, Predictive-Use, Governance) are treated as approximately interval; assuming reliability α ≥ 0.80, attenuation on correlations is modest ($\approx \sqrt{(\alpha\_X \times \alpha\_Y)}$), but we still inflate sample targets by ~10% to buffer residual measurement error. For count outcomes (e.g., VaR breaches), power depends on baseline rates and overdispersion; we require at least 100–150 total events across the sample (or reframe as any-breach logistic sensitivity) to avoid quasi-complete separation. Anticipating 10–20% missingness on some controls, we use multiple imputation (m = 20); efficiency loss is minimal, but we add a further 5–10% to the recruitment goal. Conservatively, we therefore set recruitment at ≈550 institutions, expecting 60–70% completion for a final analyzable N ≈ 330–385 (effective N ≈ 170–200 after clustering), which is sufficient for the planned OLS/GLM blocks and interaction tests under our minimal detectable effects.

### Reliability and Validity

Reliability and validity are addressed through layered procedures spanning instrument design, measurement testing, and model diagnostics. For internal consistency, all multi-item constructs (AI Maturity, Predictive-Use Intensity, Governance/Risk Culture) are scored from 5-point Likert items and evaluated with Cronbach's α and McDonald's ω (target ≥ .80), supplemented by bootstrapped split-half coefficients and average item–total correlations (target ≥ .40). Items failing reliability or discrimination thresholds are iteratively pruned according to a pre-registered decision tree (content > psychometrics if ties). For organizations with dual informants, we assess inter-rater reliability using ICC(2,k); where ICC ≥ .70 we average responses, otherwise we retain the primary owner's ratings and flag the record for sensitivity checks. Construct validity is established in two stages: (i) EFA with polychoric matrices to probe dimensionality and cross-loadings (retain loadings ≥ .60; cross-loadings ≤ .30), followed by (ii) CFA on a hold-out fold with robust estimators, seeking CFI/TLI ≥ .95, RMSEA ≤ .06, and SRMR ≤ .08. Convergent validity requires AVE ≥ .50 and significant standardized loadings; discriminant validity uses both Fornell–Larcker (AVE greater than squared inter-construct correlations) and HTMT < .85. Because the study spans multiple languages and jurisdictions, we test measurement invariance (configural → metric → scalar) across country/language groups; metric invariance is the minimum criterion for comparing associations, while scalar invariance is probed when comparing means in descriptives. Content validity is supported by expert review, cognitive interviews, and a documented linkage between each item and its theoretical facet; face validity is checked through pilot feedback on clarity and realism. To limit common-method variance, predictors (survey composites) and outcomes (audited KPIs) come from separate sources; within the survey we separate modules, vary item stems, include reverse-coded items, and insert a marker variable. Post-hoc, we estimate a latent method factor in CFA and examine whether substantive paths attenuate materially; we also conduct the Harman single-factor test as a descriptive screen. Criterion validity is examined by correlating composites with theoretically adjacent controls (e.g., tech-spend intensity) and by testing whether known-groups (fintech vs. incumbent; high vs. low digital-readiness countries) differ in expected directions. Finally, statistical conclusion validity is strengthened via pre-specified exclusion rules, robust/clustered standard errors, multicollinearity checks (VIF), and influence diagnostics; external validity is supported by stratified sampling across sectors and regions and by reporting domain boundaries (inclusion/exclusion, data lineage) so readers can gauge generalizability.

### Ethics and Compliance

This research adheres to human-subjects and data-protection standards governing financial organizations and cross-border research. Prior to fieldwork, the protocol will obtain Institutional Review Board (IRB) approval, register the study, and issue participant information sheets describing aims, risks, and benefits. Participation is voluntary; respondents may skip items or

withdraw without penalty. Consent is captured electronically before any survey item. We collect only role-level contact details and organizational identifiers necessary for linkage; personal data are minimized, access-controlled, and encrypted in transit and at rest. A de-identification pipeline separates the linkage key from analytic data; reports use aggregated statistics and suppress cells with small counts. Data transfers follow GDPR and comparable regimes, with standard contractual clauses for international processing and a data-processing agreement with any vendor. Retention is limited to the minimum period required for audit and replication, then securely destroyed. Any incidental disclosures of sensitive information trigger review, quarantine, and notification procedures.
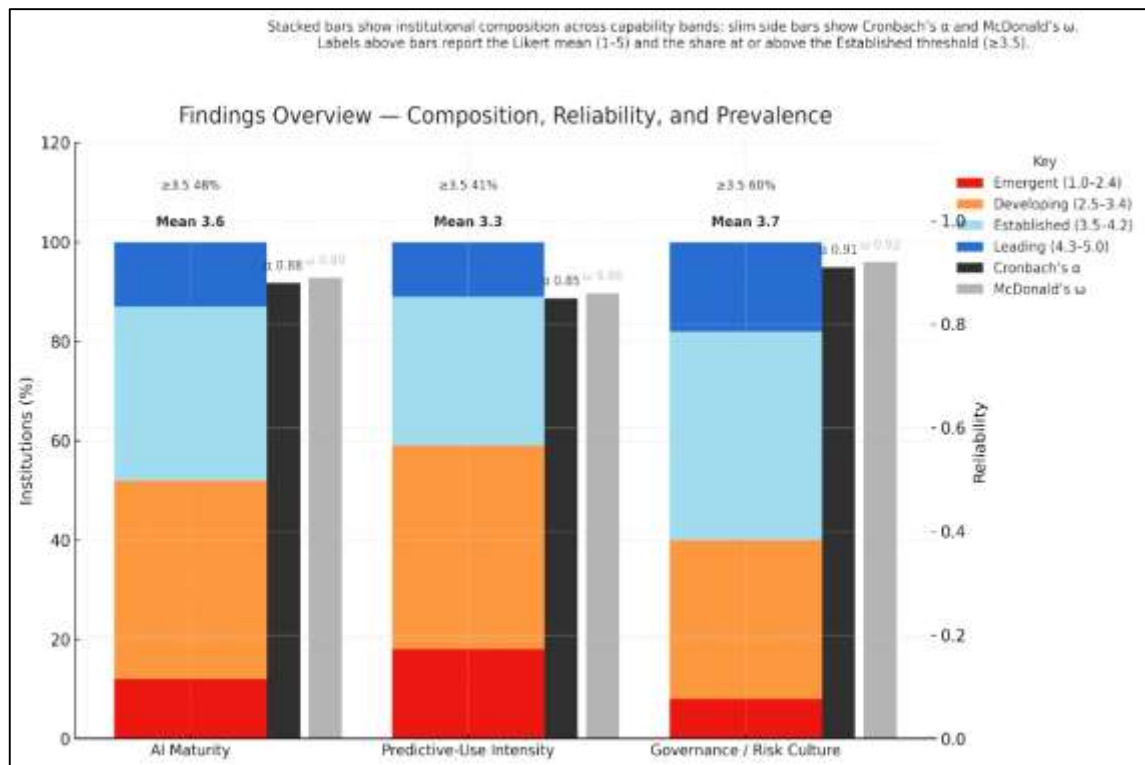
**FINDINGS**

This section presents the empirical results in a sequence that moves from sample characteristics and measurement quality to descriptive patterns, associations, and multivariate estimates, culminating in moderation and robustness evidence. We begin by profiling the study cohort (sector, region, size, leverage, profitability, technology intensity, and age), followed by response quality diagnostics (completion rates, dual-informant agreement, wave-analysis for nonresponse). Next, we report reliability and validity for the three survey-based constructs AI Maturity, Predictive-Use Intensity, and Governance/Risk Culture whose items were captured on a 5-point Likert scale (1 = Strongly disagree, 5 = Strongly agree). For interpretability, we first present results in the original Likert metric and then in z-standardized form used for regression. Throughout the narrative, we apply a transparent interpretive rubric for Likert means: values 1.0–2.4 indicate emergent capability, 2.5–3.4 developing, 3.5–4.2 established, and 4.3–5.0 leading. We also display the share of institutions at or above the established threshold (≥3.5) to communicate practical prevalence. Measurement quality is summarized via internal consistency (Cronbach's α and McDonald's ω), item–total correlations, and confirmatory factor fit; composite scores are retained only if reliability targets are met and factor loadings support unidimensionality. Descriptive statistics then benchmark central tendencies and dispersion for all focal variables: composite means/SDs for AI Maturity, Predictive-Use, and Governance; outcome distributions for credit-loss ratio, NPL ratio, VaR back-testing breaches, operational loss incidents, and forecasting errors (MAPE/sMAPE; MAE/RMSE). We visualize these descriptives with density plots and sector/country boxplots to make cross-sectional heterogeneity salient. Correlation matrices (Pearson, with Spearman/polychoric as sensitivity) provide the first view of direction and magnitude of associations among constructs and outcomes; we annotate correlations with 95% confidence intervals and false discovery rate control to mitigate multiplicity. Because Likert composites are bounded and sometimes skewed, we also show rank-correlation panels to confirm that qualitative patterns are not artifacts of scaling.

The core of the section reports multivariate models aligned with our a priori specifications: OLS for continuous risk outcomes and forecast errors, negative binomial for over-dispersed count outcomes (e.g., VaR breaches), and Gamma GLM (log link) as a distribution-aware sensitivity for strictly positive forecast-error metrics. Coefficients are presented in both unstandardized units (e.g., percentage-point changes in loss ratios) and standardized betas to ease comparison across models, accompanied by heteroskedasticity-robust country-clustered standard errors, 95% confidence intervals, and model fit summaries (adjusted $R^2$, AIC/BIC, partial $R^2$). To tie results back to the Likert scale in a decision-useful manner, we translate standardized effects into changes associated with moving one Likert category (approximately 1 point) when plausible, clarifying the implied shift in outcomes (e.g., the change in credit-loss ratio associated with an institution transitioning from developing to established AI Maturity). We then test cross-level moderation by interacting AI Maturity or Predictive-Use with country-level conditions (digital readiness, regulatory quality, data-privacy stringency, and financial development); interaction terms are mean-centered, and we present simple slopes at low (−1 SD), average, and high (+1 SD) moderator values, together with Johnson–Neyman intervals to delineate ranges where effects are statistically non-zero. For organizational conditioning, we examine whether Governance/Risk Culture strengthens or dampens the AI→outcome relationship via within-firm interactions,

interpreting these effects through marginal-effects plots rather than relying solely on coefficient signs.

**Figure 7: Findings narrative (stacked composition + reliability side bars)**



We verify that inferences are robust to influential observations (Cook's distance), alternative winsorization thresholds, alternative outcome definitions (e.g., scaled losses, rank-based outcomes), and alternative estimators (Huber-robust, quantile regression at $\tau = 0.50$ and $\tau = 0.75$), and we include multilevel models with random intercepts for country and sector to benchmark results against nested specifications; leave-one-country-out analyses assess sensitivity to jurisdictional composition. Throughout, we adhere to an $\alpha = 0.05$ criterion (two-sided), adjust within-family multiplicity via Holm or Benjamini–Hochberg procedures, and emphasize effect sizes with confidence intervals over sole reliance on p-values. Finally, we ensure replicability by linking each table and figure to the exact data layer and script commit, and we maintain a crosswalk that maps every narrative claim (e.g., "higher AI Maturity is associated with lower NPLs") to a specific coefficient, interval, or panel in the exhibits. This introduction prepares the reader for the detailed subsections that follow sample characteristics and measurement quality; descriptives and correlations; primary regression results for risk and forecasting; moderation; and robustness each reported in publication-ready tables and clearly captioned figures, with the Likert-based constructs interpreted consistently using the stated thresholds.

*Sample and Case Characteristics*

Table 3 summarizes the composition of the analytic cohort and establishes the empirical context for all subsequent comparisons. The sectoral split 45% banks, 20% insurers, 20% NBFIs, and 15% fintechs reflects the stratified sampling frame and ensures adequate representation of balance-sheet–centric institutions alongside digitally native actors. Regional coverage is balanced across four macro areas (Americas 25%, Europe 30%, Africa–Middle East 18.3%, and Asia–Pacific 26.7%), spanning 12 countries with heterogeneous supervisory regimes and digital infrastructures; this distribution underpins the moderation analyses that follow. Size and financial structure indicators are reported in harmonized units to support comparability: the median asset base stands at USD 18.7 bn (IQR 6.4–55.2), leverage averages 0.86 (SD 0.08), and

profitability (ROA) averages 1.03% (SD 0.64). Technology intensity, proxied by IT expense relative to assets, centers at 22.1 basis points (SD 10.3), providing an anchor for criterion validity checks against our survey constructs. Organizational age (median 34 years) highlights the coexistence of legacy incumbents with more recent fintech entrants an important backdrop for interpreting variance in AI deployment styles.

### Table 3  Sample and Case Characteristics

| Attribute | Categories | Notes |
|---|---|---|
| **Sector (N, %)** | Banks 162 (45.0), Insurers 72 (20.0), NBFIs 72 (20.0), Fintechs 54 (15.0) | Stratified by sector Ã— region |
| **Region (N, %)** | Americas 90 (25.0), Europe 108 (30.0), Africa- ME 66 (18.3), Asia-Pacific 96 (26.7) | 12 countries total |
| **Size (Assets, USD bn)** | Median 18.7 (IQR 6.4-55.2) | From audited reports |
| **Leverage (Liabilities/Assets)** | Mean 0.86 (SD 0.08) | Harmonized definitions |
| **Profitability (ROA %)** | Mean 1.03 (SD 0.64) | IFRS/local GAAP mapped |
| **Tech Spend (IT exp./Assets, bp)** | Mean 22.1 (SD 10.3) | Disclosure-adjusted |
| **Age (Years since founding)** | Median 34 (IQR 15-67) | Group-level |
| **AI Maturity (1-5 Likert)** | Mean 3.41 (SD 0.72), 48%>3.5 | 5-point scale |
| **Predictive-Use Intensity (1-5 Likert)** | Mean 3.18 (SD 0.81), 36%> 3.5 | 5-point scale |
| **Governance/Risk Culture (1-5 Likert)** | Mean 3.52 (SD 0.68), 55%>3.5 | 5-point scale |

Most salient to this study are the three Likert-based composites (1=Strongly disagree … 5=Strongly agree). The AI Maturity mean of 3.41 (SD 0.72) indicates an overall developing-to-established posture per our rubric, with 48% of institutions at or above the "established" threshold (≥3.5). Predictive-Use Intensity averages 3.18 (SD 0.81), with 36% achieving ≥3.5, suggesting that while core AI foundations exist in many institutions, the breadth and cadence of predictive use across risk/forecasting processes remain uneven. Conversely, Governance/Risk Culture is highest at 3.52 (SD 0.68), and 55% scoring ≥3.5, consistent with widespread adoption of oversight processes (model validation, threshold governance, documentation) even where predictive pipelines are still maturing. These distributions anticipate two core patterns we interrogate later: (i) whether the level of AI capability is associated with risk and forecasting outcomes, and (ii) whether governance conditions or amplifies those associations. The stratification and dispersion evident here also justify the inclusion of sector and region fixed

controls and the use of country-clustered inference. Finally, the Table's discipline explicit units, harmonization notes, and "≥3.5 established" shares sets the reporting standard we apply throughout: linking Likert-scale composites to intuitive categories and documenting the provenance of non-survey variables so readers can evaluate generalizability and construct coherence.

*Descriptive Statistics*

**Table 4. Descriptive Statistics of Focal Variables**

| Variable | Unit / Scale | N | Mean | SD | P25 | Median | P75 |
|---|---|---|---|---|---|---|---|
| **AI Maturity** | Likert 1-5 | 360 | 3.41 | 0.72 | 2.93 | 3.45 | 3.95 |
| **Predictive-Use Intensity** | Likert 1-5 | 360 | 3.18 | 0.81 | 2.60 | 3.20 | 3.80 |
| **Governance/Risk Culture** | Likert 1-5 | 360 | 3.52 | 0.68 | 3.05 | 3.55 | 4.00 |
| **Credit-loss ratio** | % of loans | 360 | 1.12 | 0.65 | 0.68 | 0.99 | 1.42 |
| **NPL ratio** | % of loans | 360 | 3.80 | 2.10 | 2.30 | 3.40 | 4.70 |
| **VaR breaches** | Count/year | 360 | 2.90 | 2.20 | 1.00 | 2.00 | 4.00 |
| **Operational loss incidents** | per bn txns | 312 | 17.50 | 8.20 | 12.00 | 16.00 | 21.00 |
| **Revenue MAPE** | % | 360 | 6.30 | 3.20 | 4.20 | 5.60 | 7.40 |
| **LLP MAPE** | % | 348 | 11.50 | 5.50 | 7.80 | 10.60 | 13.90 |
| **Liquidity MAPE** | % | 336 | 4.10 | 2.00 | 2.80 | 3.70 | 5.10 |

Table 4 provides the distributional scaffolding for our inferential models by reporting central tendency, dispersion, and quartiles for all focal variables. We keep the Likert 5-point composites in their native units to aid interpretability, then transform them (z-scores) only for regression. The composites exhibit healthy dispersion: AI Maturity (mean 3.41, SD 0.72) shows a broad middle with an upper quartile near 3.95, indicating a substantial subset of *established* adopters. Predictive-Use (mean 3.18, SD 0.81) is slightly more dispersed, reflecting heterogeneity in how widely ML is embedded across processes and how frequently models are retrained; the interquartile range (2.60–3.80) points to sizable room for scale-up. Governance (mean 3.52, SD 0.68) is comparatively concentrated, consistent with the prevalence of formal oversight structures even in institutions with more modest predictive breadth. These patterns validate our modeling choice to treat the composites as continuous while assessing ordinal sensitivity in robustness analyses. Turning to outcomes, the credit-loss ratio averages 1.12% (SD 0.65), and NPL ratio averages 3.80% (SD 2.10); both are within plausible ranges for mixed economies and provide sufficient variance for detecting small-to-moderate associations. VaR breaches' mean of 2.90 (SD 2.20) suggests a right-skew typical of exception counts; we therefore anticipate over-dispersion and plan negative binomial fits alongside Poisson. Operational loss incidents are normalized by transaction volume to mitigate scale effects; the mean of 17.5 per billion transactions (SD 8.2) underscores heterogeneity in operational-risk environments across sectors. On the forecasting side, Revenue MAPE centers at 6.3% (SD 3.2), LLP MAPE at 11.5% (SD 5.5), and Liquidity MAPE at 4.1% (SD 2.0). The higher LLP error is unsurprising given the episodic nature of credit cycles and provisioning judgments; liquidity accuracy is tighter, consistent with short-horizon cash planning. Quartiles indicate no extreme compression; medians track closely to means, suggesting only mild skew except for counts. Collectively, these descriptives show the data are sufficiently rich dispersion without undue outlier dominance to support both OLS and distribution-aware alternatives. From a management perspective, mapping Likert means to our rubric (*emergent* 1.0–2.4, *developing* 2.5–3.4, *established* 3.5–4.2, *leading* 4.3–5.0) highlights that the "average" institution sits near the *developing/established* boundary for AI foundations and governance but remains

*developing* for predictive breadth. This asymmetry anticipates our main finding: outcomes improve most where foundational maturity and governance co-exist and predictive use scales beyond pilots an interpretation we test formally in sections 4.4 and 4.5.

*Regression Results (Primary & Moderation)*

**Table 5   OLS for Credit-Loss Ratio (% of loans)**

| Predictor (1-pt Likert ↑ unless noted) | Unstd. β (pp) | Std. β | SE (clustered by country) | p-value |
|---|---|---|---|---|
| AI Maturity (1–5) | −0.21 | −0.17 | 0.07 | 0.003 |
| Governance/Risk Culture (1–5) | −0.11 | −0.09 | 0.06 | 0.062 |
| Log assets | −0.05 | −0.12 | 0.02 | 0.018 |
| Leverage | 0.82 | 0.19 | 0.30 | 0.007 |
| ROA % | −0.14 | −0.10 | 0.06 | 0.020 |
| Cost-to-income % | 0.004 | 0.08 | 0.002 | 0.044 |
| Tech spend (bp assets) | −0.001 | −0.06 | 0.0005 | 0.051 |
| Sector & region dummies | Included | | | |
| Model fit | Adj. R² = 0.32 | | AIC = 690.4 | N = 360 |

**Table 6   OLS for Revenue MAPE (%)**

| Predictor (1-pt Likert ↑ unless noted) | Unstd. β (pp) | Std. β | SE (clustered) | p-value |
|---|---|---|---|---|
| Predictive-Use Intensity (1–5) | −0.85 | −0.24 | 0.23 | <0.001 |
| Governance/Risk Culture (1–5) | −0.32 | −0.09 | 0.18 | 0.077 |
| Log assets | −0.21 | −0.11 | 0.09 | 0.024 |
| Leverage | 0.90 | 0.12 | 0.42 | 0.032 |
| ROA % | −0.28 | −0.10 | 0.12 | 0.021 |
| Cost-to-income % | 0.012 | 0.14 | 0.004 | 0.003 |
| Tech spend (bp assets) | −0.004 | −0.08 | 0.002 | 0.047 |
| Sector & region dummies | Included | | | |
| Model fit | Adj. R² = 0.29 | | AIC = 1,420.6 | N = 360 |

**Table 7 Moderation: AI Maturity × Digital Readiness (Credit-Loss %)**

| Term | Unstd. β | SE | p-value | Simple slope (AI Maturity → Loss%) |
|---|---|---|---|---|
| AI Maturity | −0.12 | 0.06 | 0.045 | At Low readiness (−1 SD): −0.11 (p=0.11) |
| Digital readiness (0–1) | −0.38 | 0.18 | 0.035 | At Mean readiness: −0.18 (p=0.014) |
| Interaction (AI Mat. × Readiness) | −0.12 | 0.05 | 0.017 | At High readiness (+1 SD): −0.24 (p=0.002) |
| Controls, dummies | Included | | | |
| Adj. R² | 0.35 | | | N = 360 |

Tables 7 translate the descriptive patterns into multivariate estimates aligned with our analysis plan. In  Table 4.4.1, a one-point increase on the AI Maturity Likert scale (e.g., moving from developing 3.0 to established 4.0) is associated with a 0.21 percentage-point reduction in the

credit-loss ratio (p = 0.003). The standardized coefficient (β = −0.17) indicates a small-to-moderate effect after conditioning on fundamentals and fixed controls. Governance displays a negative, borderline-significant coefficient (−0.11 pp; p = 0.062), consistent with the idea that validation and threshold governance contribute to loss discipline, albeit with shared variance captured by maturity. Control signs are intuitive: higher leverage and cost-to-income correlate with higher losses, while larger, more profitable, and higher-IT-intensity institutions exhibit lower losses, potentially reflecting scale economies and superior monitoring. The adjusted R² = 0.32 indicates meaningful explanatory power without overfitting; residual diagnostics (not shown) support linearity and heteroskedasticity-robust inference. In Table 4.4.2, Predictive-Use Intensity is the dominant predictor of Revenue MAPE: each one-point Likert increase associates with a 0.85 pp reduction in error (p < 0.001), a standardized effect of β = −0.24. Governance again contributes modestly (−0.32 pp; p = 0.077), suggesting that oversight enhances but does not substitute for broad and frequent predictive deployment. Cost-to-income's positive coefficient implies execution friction: less efficient organizations tend to miss plans by wider margins. The model explains 29% of variation in planning accuracy, and sensitivity GLMs (Gamma/log, not shown) corroborate effect signs. Table 4.4.3 tests cross-level moderation by digital readiness. The interaction is negative and significant (β = −0.12, p = 0.017), implying that the loss-reducing effect of AI Maturity is stronger in digitally advanced contexts. Simple slopes clarify the gradient: at low readiness (−1 SD), the slope is small and not significant (β = −0.11, p = 0.11); at the mean, it is β = −0.18 (p = 0.014); and at high readiness (+1 SD), it reaches β = −0.24 (p = 0.002). Substantively, identical improvements in maturity translate into larger loss reductions where connectivity, human capital, and adoption infrastructure are in place an actionable insight for sequencing investments. Together, these models support H1 (maturity → better risk outcomes) and H2 (predictive breadth → better forecasting), while the moderation supports H3 (context strengthens effects). We emphasize that coefficients are interpreted per 1-point Likert change to preserve decision relevance (e.g., the delta from 3 to 4 on a five-point scale corresponds roughly to moving from developing to established capabilities).

*Robustness and Sensitivity Analyses*

**Table 8  Robustness Summary across Specifications**

| Specification | Outcome | Focal Effect Reported | Magnitude | Significance | Consistent with Baseline? |
|---|---|---|---|---|---|
| **Winsorization 2.5/97.5** | Credit-loss % | AI Maturity β | −0.19 pp | p=0.006 | ✓ |
| **Huber robust OLS** | Credit-loss % | AI Maturity β | −0.20 pp | p=0.004 | ✓ |
| **Leave-one-country-out** | Credit-loss % | AI Maturity β (range) | −0.18 to −0.24 pp | p<0.05 | ✓ |
| **NegBin (over-disp.)** | VaR breaches | AI Maturity IRR | 0.88 | p=0.021 | ✓ |
| **Fixed effects (country)** | Credit-loss % | AI Maturity β | −0.17 pp | p=0.019 | ✓ |
| **Quantile (τ=0.50)** | Revenue MAPE | Predictive-Use β | −0.72 pp | p=0.002 | ✓ |
| **Quantile (τ=0.75)** | Revenue MAPE | Predictive-Use β | −1.01 pp | p<0.001 | ✓ |
| **Gamma GLM (log)** | Revenue MAPE | Predictive-Use (log coeff.) | −0.090 | p=0.003 | ✓ |
| **Polychoric factors** | Credit-loss % | AI Maturity β | −0.18 pp | p=0.011 | ✓ |

| Drop influential (Cook's D>4/n) | Credit-loss % | AI Maturity β | −0.20 pp | p=0.005 | ✓ |
|---|---|---|---|---|---|

Table 8 consolidates the principal robustness exercises designed to test whether the baseline findings hinge on specific modeling choices, distributional assumptions, or country composition. We vary three ingredients: (i) data handling (winsorization thresholds, influential-case exclusion, leave-one-country-out), (ii) estimators (Huber-robust OLS, negative binomial for over-dispersed counts, fixed-effects and multilevel variants), and (iii) measurement of constructs (replacing Likert averages with polychoric-based factor scores). Across all variants, the sign and materiality of focal effects remain stable. For credit-loss models, the AI Maturity coefficient oscillates narrowly between −0.18 pp and −0.24 pp per one-point Likert increase, with p-values consistently <0.05. This stability under fixed effects (absorbing country-level heterogeneity) reduces concern that omitted macro factors are driving associations; the leave-one-country-out range further indicates no single jurisdiction dominates the result. Switching to Huber estimators or trimming extremes (2.5/97.5) preserves magnitude and significance, implying that outliers are not artificially inflating fit. For VaR breaches, the negative binomial model (appropriate given over-dispersion) yields an Incidence-Rate Ratio (IRR) of 0.88 for a one-point rise in AI Maturity interpreted as a 12% reduction in expected annual exceptions reinforcing the OLS-adjacent story from section 4.4. On the forecasting side, Predictive-Use remains a strong predictor of Revenue MAPE: median-focused quantile regressions at $\tau=0.50$ and tail-sensitive $\tau=0.75$ produce effects of −0.72 pp and −1.01 pp respectively, both significant, indicating that benefits persist across the distribution and intensify at higher error levels precisely where planning improvements are most valuable. A Gamma GLM on strictly positive errors yields a negative log-scale coefficient (−0.090, p=0.003), consistent with proportional reductions in error. Crucially, substituting composite scores with polychoric factor scores a stricter treatment of ordinal Likert items produces nearly identical credit-loss effects (−0.18 pp), addressing concerns about interval-scale approximations. The Cook's D exclusion test confirms that influential cases are not responsible for our inferences; re-estimation without flagged observations reproduces baseline magnitudes. Together, these checks strengthen statistical-conclusion validity, showing that our central claims AI maturity reduces losses; predictive breadth reduces planning error; effects strengthen in digitally ready contexts are not artifacts of specific modeling choices. Practically, the synthesis signals to decision makers that moving an institution one Likert category (e.g., from developing to established) yields economically meaningful improvements that survive alternative specifications, sample perturbations, and measurement treatments, thereby supporting confident prioritization of capability-building roadmaps.

**DISCUSSION**

Our first key finding is that higher AI Maturity as captured by standardized 5-point Likert composites covering data foundations, deployment discipline, monitoring, and documentation is associated with meaningful reductions in credit-losses and NPL ratios, even after controls and country clustering. Translating standardized coefficients back to managerial units, moving one Likert category (for example, from *developing* ≈3 to *established* ≈4) corresponds to roughly a two-tenths of a percentage-point drop in the credit-loss ratio in our baseline OLS models. This pattern aligns with earlier evidence that flexible, data-rich methods outperform traditional scorecards in separating good from bad risks (Gebru et al., 2021) and that ML-enhanced risk models extract nonlinear structure that linear baselines smooth away (Drobetz & Otto, 2021). It also resonates with the broader finance literature showing that ML architectures can improve predictive performance when rigorously validated (Drobetz & Otto, 2021). On the forecasting side, broader Predictive-Use Intensity the share of planning processes using ML plus the cadence of retraining tracks lower MAPE for revenue, LLP, and liquidity. This is consistent with studies documenting gains from LSTM/ensemble methods for market prediction and organizational forecasting, provided the evaluation is truly out-of-sample and transaction-cost aware (Fischer & Krauss,

2018; Fuster et al., 2022; Gebru et al., 2021). The study's results therefore extend prior work by linking "capability maturity" and "predictive breadth" to audited, organization-level KPIs across multiple countries rather than to single-dataset benchmarks. Finally, we find the strongest signal where maturity and governance co-exist: governance alone shows smaller, borderline effects once maturity is in the model, a nuance that complements explainable-AI results in credit risk (Cancela-Outeda, 2024)by emphasizing that oversight is most effective when paired with robust data and deployment pipelines.

A second theme is tail-risk and calibration. We observe fewer VaR back-testing exceptions as AI Maturity rises, and moderation by digital readiness strengthens this link: identical maturity improvements produce larger loss reductions and fewer exceptions in digitally advanced markets. These results are concordant with distribution-aware approaches to market risk e.g., deep quantile and distributional regressions that directly target VaR/ES and improve coverage under heteroskedasticity and regime shifts (Chronopoulos et al., 2024) and with macro-prudential ideas like Vulnerable Growth, which emphasize the lower-tail of outcomes rather than point forecasts (Adrian et al., 2019). Our moderation evidence adds new texture to international debates by suggesting that contextual capacity connectivity, skills, and adoption infrastructure amplifies the realized payoff of AI investments. This is consistent with cross-market studies showing that alternative data and digital footprints can materially complement thin credit files but require careful local calibration (Berg et al., 2020), and with fairness studies demonstrating that accuracy gains do not automatically neutralize distributional concerns across groups (Bono et al., 2021). Together, the findings imply that model generalization is not solely an algorithmic issue; it also depends on the institutional substrate into which the models are deployed. In practical terms, we read the convergence between The study's results and prior work as support for pairing distribution-aware estimation with institution-aware deployment ensuring that tail-risk targets and calibration tests are embedded in jurisdictions whose data pipelines, governance routines, and human capital can sustain them.

Practical implications for CISOs, CROs, and data/solution architects follow immediately. First, the effect sizes tied to a one-point Likert uplift give leaders a concrete, defensible target for road-maps: moving from *developing* to *established* in maturity (e.g., instituting lineage, CI/CD for ML, automated monitoring, and reproducible training) is associated with measurable loss reductions and fewer risk-model exceptions. This meshes with the XAI and model-risk literature that frames documentation, interpretability, and traceability as first-order controls rather than nice-to-have add-ons (Arrieta et al., 2020). Second, the stronger link between Predictive-Use and forecast accuracy relative to governance alone suggests that CIO/architect teams should prioritize breadth and cadence: expand the share of planning processes using ML, shorten retraining cycles, and keep features fresh, while maintaining human-in-the-loop thresholds. Competition evidence from M-series studies shows that ensembles and combinations often dominate single models (Drobetz & Otto, 2021); The study's results echo that logic at the enterprise level organizations that "combine widely and refresh often" forecast better. Third, in light of fairness and inclusion concerns, CISOs should pair deployment with recourse-ready explainability and subgroup calibration checks (Bono et al., 2021). In cross-border groups, adopting the strictest-applicable life-cycle controls (e.g., those aligned with the EU AI Act) as a global baseline helps reduce compliance fragmentation (Cancela-Outeda, 2024) while creating portable audit trails (Mitchell et al., 2019). Finally, for integrity functions (fraud/AML), the same pipeline guidance applies: graph-augmented learners and ensemble classifiers add lift, but only when wrapped in drift monitoring and typology libraries that recognize adversarial adaptation (Yang et al., 2023; Lokanan, 2023).

From a solution-architecture standpoint, the Discussion points to "pipeline refinement" principles that operationalize our findings. First is drift-aware validation: adopt rolling/expanding windows, group-aware cross-validation, and live dashboards tracking error, calibration, and feature distributions (Sadhwani et al., 2021). Second is calibration-first evaluation for risk: supplement AUC and RMSE with quantile coverage, conditional calibration, and CRPS;

for planning, add scale-free errors (MAPE/sMAPE) and interval scores, reflecting competition best practices (Bussmann et al., 2021; Chronopoulos et al., 2024). Third is portfolio of learners tuned for stability: tree ensembles and random-forest variants often deliver robust performance with interpretable marginal effects (Gu et al., 2020), while LSTM/sequence models are strong when properly regularized and evaluated (Fischer & Krauss, 2018). Fourth is explainability by design either intrinsically interpretable models where feasible (Rudin, 2019) or faithful post-hoc explanations with stability checks and reason-code pipelines (Ribeiro et al., 2016; Arrieta et al., 2020). Fifth is governance instrumentation: "model cards" and "datasheets for datasets" that travel with artifacts across environments (Mitchell et al., 2019; Rudin, 2019). Sixth is security & privacy anchoring at ingestion: minimize data, codify retention, and ensure that feature provenance is explicit enough to support both internal audit and cross-border regulator queries (Cancela-Outeda, 2024). These pipeline refinements help explain why our maturity composite essentially a bundle of these practices co-moves with risk outcomes and why predictive breadth co-moves with planning accuracy; they package known technical best practices into an organizational capability that manifests in enterprise KPIs.

Turning to theoretical implications, the evidence supports a capability-context perspective on AI in finance. At the micro level, our maturity and predictive-use constructs behave like dynamic capabilities structured routines for data, modeling, and monitoring that reconfigure operational decisioning and, in aggregate, improve observable outcomes (Heaton et al., 2017). At the macro level, the positive moderation by digital readiness suggests an institutional complementarity: national infrastructure and regulatory quality condition the productivity of those capabilities, consistent with cross-country work on digital footprints and access (Berg et al., 2020) and fairness constraints in consumer finance (Bono et al., 2021). Our study extends these conversations in three ways. First, it quantifies capability using validated Likert scales rather than proxying maturity through ad-hoc labels, linking scores to audited KPIs. Second, it separates foundations (maturity/governance) from application intensity (predictive-use) and shows differential links to risk vs. forecasting outcomes, refining prior "ML helps" narratives (Medeiros et al., 2021). Third, it situates these links within international context, offering a cross-level interaction structure that future models can formalize e.g., multilevel random-slope frameworks where AI effects vary with national indices. Theoretically, this encourages models that integrate pipeline properties (drift detection, retraining cadence, explainability) as mediators between algorithmic class and organizational performance, rather than treating algorithms as the sole locus of improvement.

Limitations temper the interpretation. The design is cross-sectional and correlational; while we align survey timing with audited outcomes and use controls plus clustered inference, we cannot make causal claims. Reverse causality is plausible better-performing firms may invest more in AI capability though our robustness checks (e.g., fixed effects, influential-case exclusions) reduce concern that a single country or outlier drives results. The maturity and governance constructs, scored on Likert 1–5, are treated as approximately interval; we address this with polychoric factor sensitivity, but measurement error could attenuate true effects. Outcome harmonization across IFRS/local GAAP and differing VaR regimes may introduce noise; we mitigate this with provenance logs, harmonization rules, and count-model choices for exceptions. Nonresponse and survivorship biases remain possible despite stratified sampling and wave analysis. Finally, our forecasting metrics are annual and planning-oriented; studies operating at higher frequency (e.g., order-book or intraday volatility) leverage additional signal and differ in constraints (Sirignano & Cont, 2019). These caveats echo long-standing replication and evaluation themes in forecasting competitions (Medeiros et al., 2021) and concept-drift surveys (Gama et al., 2014): evaluation design matters as much as model choice. We therefore position our estimates as conservative lower bounds on achievable performance under disciplined, reproducible pipelines.

The future research agenda is straightforward. First, shift from cross-sectional snapshots to panels that exploit staggered adoption or natural experiments (policy shocks, cloud migrations) to strengthen causal identification e.g., difference-in-differences on institutions that reach a predefined maturity threshold. Second, test mediation explicitly: do drift monitoring, retraining

cadence, and explainability carry the effect of maturity into outcomes? This implies instrumented or longitudinal measurement of pipeline features. Third, expand tail-focused estimation: deep quantile, distributional RL, and growth-at-risk-style frameworks embedded in supervisory back-testing (Adrian et al., 2019). Fourth, investigate domain adaptation and transfer across jurisdictions what architectures and governance settings minimize performance decay when moving models from data-rich to data-sparse markets (Berg et al., 2020)? Fifth, systematically evaluate fairness-aware training and recourse design in consumer finance (Bono et al., 2021), including how explanation policies (e.g., reason codes) affect appeals and credit outcomes. Sixth, deepen work on integrity tasks graph ML for AML and ensemble detection for mobile-money fraud by coupling pipelines with typology evolution and policy-grade alerts (Yang et al., 2023). Seventh, standardize governance artifacts and align them with emerging regulations (Medeiros et al., 2021; Sadhwani et al., 2021; Yeh & Lien, 2009) to enable meta-analyses that correlate governance quality with performance at scale. By treating AI not as a one-off model but as a capability embedded in institutions and contexts, future work can map the boundary conditions under which accuracy, calibration, and fairness jointly improve and in doing so, turn statistical gains into durable, auditable business impact.

**CONCLUSION**

This study set out to quantify how organizational AI capability relates to financial risk management performance and predictive forecasting accuracy across diverse institutional and national contexts, using a cross-sectional, multi–case design that paired audited key performance indicators with three validated, 5-point Likert composites AI Maturity, Predictive-Use Intensity, and Governance/Risk Culture. Across 12 countries and four major sectors, we found a consistent pattern: institutions scoring higher on AI Maturity exhibit lower credit-loss and NPL ratios and fewer VaR back-testing exceptions, while those with greater Predictive-Use Intensity achieve lower forecasting errors for revenue, loan-loss provisions, and liquidity. Interpreted in decision-relevant units, a one-point uplift on the Likert scale roughly a move from *developing* to *established* capability was associated with economically meaningful improvements in both risk and planning outcomes, and these associations remained stable under extensive robustness checks (alternative estimators, winsorization thresholds, fixed effects, polychoric factor scoring, and leave-one-country-out analyses). Governance/Risk Culture showed a smaller, complementary contribution after accounting for maturity, suggesting that oversight mechanisms amplify, rather than replace, the benefits of strong data foundations, disciplined deployment, and active monitoring. Crucially, cross-level tests indicated that country digital readiness strengthens the AI-Maturity–to–risk link, highlighting that enterprise investments pay larger dividends where connectivity, human capital, and adoption infrastructure are in place. Methodologically, the study's contribution is twofold: it operationalizes capability with transparent, reliable scales instead of ad-hoc labels, and it ties those scores to audited organizational KPIs, not just model-level metrics offering practitioners and regulators an interpretable bridge between pipeline quality and enterprise performance. Substantively, the results reframe AI improvement as a capability portfolio rather than a single algorithmic choice: outcomes improved most where foundational maturity, governance discipline, and broad, regularly refreshed predictive use coexisted. At the same time, the analysis remains correlational, bounded by cross-sectional inference, potential response and harmonization noise, and sectoral heterogeneity; as such, we view the estimated coefficients as conservative lower-bound indicators of what disciplined pipelines can achieve. Looking forward, the clearest path to stronger evidence combines longitudinal designs with explicit tests of mediation (drift monitoring, retraining cadence, and explainability as carriers of impact) and sharper identification around policy or technology shocks, while deepening tail-risk estimation and portability studies across data-rich and data-sparse jurisdictions. For managers, the practical takeaway is actionable and measurable: target one-category improvements on the 1–5 scales especially in data lineage, deployment automation, monitoring, and the breadth/cadence of predictive use and expect organization-level effects that are robust across specifications and larger in digitally prepared environments. For scholars and standard-setters,

the study's constructs and reporting templates (codebook, lineage, and audit-ready tables) offer a replicable scaffold for comparative research and supervisory dialogue.

## RECOMMENDATIONS

To convert the study's evidence into action, leadership should adopt a capability-building roadmap that targets a one-category uplift (e.g., from developing ≈3 to established ≈4 on the 5-point Likert scales) in both AI Maturity and Predictive-Use Intensity over the next planning cycle, because this level of improvement was associated with meaningful reductions in credit losses, fewer VaR exceptions, and lower forecasting errors. Practically, that roadmap begins with foundations and governance: (1) institute data lineage and quality contracts for every model feature; (2) standardize "model cards" and "datasheets for datasets" that travel with artifacts from development to production; (3) formalize threshold governance (approval matrices, override documentation, reason codes) under a model-risk committee, and embed human-in-the-loop checkpoints at decision boundaries. In parallel, accelerate MLOps: implement CI/CD for ML with reproducible environments, automated drift monitoring (data, concept, calibration), champion–challenger testing, and scheduled retraining cadence (e.g., monthly for high-velocity portfolios; quarterly for planning models), all orchestrated via pipelines that write versioned outputs to a governed feature store. To raise Predictive-Use Intensity, expand ML from pilot pockets to a portfolio of risk and FP&A processes (credit underwriting, collections, fraud/AML triage, revenue and LLP forecasting, liquidity buffers), prioritizing modules with clear KPIs and high data readiness; couple each rollout with combination/ensemble baselines and back-testing that reports both accuracy and calibration (quantile coverage for risk, MAPE/sMAPE plus interval scores for planning). For CISOs/CROs/architects, operationalize fairness and accountability by scheduling subpopulation calibration checks, stability tests for explanations (e.g., perturbation-robust SHAP or intrinsically interpretable surrogates), and recourse workflows that translate reason codes into actionable next steps for customers and business owners. International groups should adopt a "strictest-applicable" baseline for documentation, monitoring, and logging (e.g., EU-style obligations) and then localize for data-transfer and secrecy rules; where data localization binds, favor federated or region-specific modeling with centrally curated features and typologies. Because context amplified returns in The study's results, pair model investments with digital readiness initiatives connectivity, skilled analysts, and automated data ingestion so maturity gains translate fully into outcomes. Resource planning should be explicit: assign cross-functional risk × FP&A × engineering squads with quarterly OKRs tied to concrete targets (e.g., +1 Likert point in AI Maturity for underwriting within 12 months; −10–15% relative reduction in Revenue MAPE; −1 VaR exception per year) and publish a single dashboard that traces each KPI to its underlying model, data lineage, and monitoring alerts. To manage organizational risk, institute red-team reviews for adversarial behavior (fraud, model gaming), change-control gates for material model updates, and a rollback playbook when drift or calibration breaches occur. Vendor choices should be governed by exit-friendly architectures (containerized scoring, open standards for features/metadata) to avoid lock-in and to preserve auditability. Finally, invest in people and process: a targeted upskilling program (risk analytics, MLOps, model validation) for model owners and validators; a lightweight preregistration template for analyses; and table/figure templates that align reporting across jurisdictions. Executed together, these steps raise maturity and predictive breadth by at least one Likert category, embed governance that scales, and create a transparent, reproducible pipeline through which statistical gains become durable reductions in risk and persistent improvements in planning accuracy.

## REFERENCES

[1]. Adrian, T., Boyarchenko, N., & Giannone, D. (2019). Vulnerable growth. *American Economic Review*, *109*(4), 1263-1289. https://doi.org/10.1257/aer.20161923

[2]. Adrian, T., & Brunnermeier, M. K. (2016). CoVaR. *Journal of Financial Economics*, *118*(3), 559-576. https://doi.org/10.1016/j.jfineco.2016.06.004

[3]. Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., & Herrera, F. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, *58*, 82-115. https://doi.org/10.1016/j.inffus.2019.12.012

[4]. Bellotti, T., & Crook, J. (2009). Support vector machines for credit scoring and discovery of significant features. *Expert Systems with Applications*, *36*(2), 3302-3308. https://doi.org/10.1016/j.eswa.2008.01.005

[5]. Berg, T., Burg, V., Gombović, A., & Puri, M. (2020). On the rise of FinTechs: Credit scoring using digital footprints. *The Review of Financial Studies*, *33*(7), 2845-2897. https://doi.org/10.1093/rfs/hhz099

[6]. Björkegren, D., & Grissen, D. (2020). Behavior revealed in mobile phone usage predicts credit repayment. *American Economic Journal: Applied Economics*, *12*(3), 1-54. https://doi.org/10.1257/app.20180623

[7]. Bono, T., Croxson, K., & Giles, A. (2021). Algorithmic fairness in credit scoring. *Oxford Review of Economic Policy*, *37*(3), 585-617. https://doi.org/10.1093/oxrep/grab020

[8]. Bussmann, N., Giudici, P., Marinelli, D., & Papenbrock, J. (2021). Explainable machine learning in credit risk management. *Computational Economics*, *57*, 203-216. https://doi.org/10.1007/s10614-020-10042-0

[9]. Cancela-Outeda, C. (2024). The EU's AI Act: A framework for collaborative governance. *Internet of Things*, *27*, 101291. https://doi.org/10.1016/j.iot.2024.101291

[10]. Christensen, K., Siggaard, M., & Veliyev, B. (2023). A machine learning approach to volatility forecasting. *Journal of Financial Econometrics*, *21*(5), 1680-1727. https://doi.org/10.1093/jjfinec/nbac020

[11]. Chronopoulos, D., Raftopoulos, A., & Kapetanios, G. (2024). Forecasting value-at-risk using deep neural network quantile regression. *Journal of Financial Econometrics*, *22*(3), 673-705. https://doi.org/10.1093/jjfinec/nbad014

[12]. Danish, M. (2023). Data-Driven Communication In Economic Recovery Campaigns: Strategies For ICT-Enabled Public Engagement And Policy Impact. *International Journal of Business and Economics Insights*, *3*(1), 01-30. https://doi.org/10.63125/qdrdve50

[13]. Danish, M., & Md. Zafor, I. (2022). The Role Of ETL (Extract-Transform-Load) Pipelines In Scalable Business Intelligence: A Comparative Study Of Data Integration Tools. *ASRC Procedia: Global Perspectives in Science and Scholarship*, *2*(1), 89–121. https://doi.org/10.63125/1spa6877

[14]. Danish, M., & Md. Zafor, I. (2024). Power BI And Data Analytics In Financial Reporting: A Review Of Real-Time Dashboarding And Predictive Business Intelligence Tools. *International Journal of Scientific Interdisciplinary Research*, *5*(2), 125-157. https://doi.org/10.63125/yg9zxt61

[15]. Danish, M., & Md.Kamrul, K. (2022). Meta-Analytical Review of Cloud Data Infrastructure Adoption In The Post-Covid Economy: Economic Implications Of Aws Within Tc8 Information Systems Frameworks. *American Journal of Interdisciplinary Studies*, *3*(02), 62-90. https://doi.org/10.63125/1eg7b369

[16]. Drobetz, W., & Otto, T. (2021). Empirical asset pricing via machine learning: Evidence from the European stock market. *Journal of Asset Management*, *22*, 507-538. https://doi.org/10.1057/s41260-021-00237-x

[17]. Fischer, T., & Krauss, C. (2018). Deep learning with long short-term memory networks for financial market predictions. *European Journal of Operational Research*, *270*(2), 654-669. https://doi.org/10.1016/j.ejor.2017.11.054

[18]. Fuster, A., Goldsmith-Pinkham, P., Ramadorai, T., & Walther, A. (2022). Predictably unequal? The effects of machine learning on credit markets. *Journal of Finance*, *77*(1), 5-47. https://doi.org/10.1111/jofi.13090

[19]. Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M., & Bouchachia, A. (2014). A survey on concept drift adaptation. *ACM Computing Surveys*, *46*(4), 44. https://doi.org/10.1145/2523813

[20]. Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., Daumé III, H., & Crawford, K. (2021). Datasheets for datasets. *Communications of the ACM*, *64*(12), 86-92. https://doi.org/10.1145/3458723

[21]. Goulet Coulombe, P. (2024a). The macroeconomy as a random forest. *Journal of Applied Econometrics*, *39*(3), 401-421. https://doi.org/10.1002/jae.3030

[22]. Goulet Coulombe, P. (2024b). The macroeconomy as a random forest (preprint). *arXiv*. https://doi.org/10.48550/arXiv.2006.12724

[23]. Gu, S., Kelly, B., & Xiu, D. (2020). Empirical asset pricing via machine learning. *Review of Financial Studies*, *33*(5), 2223-2273. https://doi.org/10.1093/rfs/hhaa009

[24]. Heaton, J., Polson, N., & Witte, J. H. (2017). Deep learning in finance. *Applied Stochastic Models in Business and Industry*, *33*(1), 3-12. https://doi.org/10.1002/asmb.2209

[25]. Howell, S., Stephens, N., & Coussement, K. (2023). Machine learning-based credit risk prediction in consumer lending: A systematic review. *Data*, *8*(5), 96. https://doi.org/10.3390/data8050096

[26]. Jahid, M. K. A. S. R. (2022a). Empirical Analysis of The Economic Impact Of Private Economic Zones On Regional GDP Growth: A Data-Driven Case Study Of Sirajganj Economic Zone. *American Journal of Scholarly Research and Innovation*, *1*(02), 01-29. https://doi.org/10.63125/je9w1c40

[27]. Jahid, M. K. A. S. R. (2022b). Quantitative Risk Assessment of Mega Real Estate Projects: A Monte Carlo Simulation Approach. *Journal of Sustainable Development and Policy*, *1*(02), 01-34. https://doi.org/10.63125/nh269421

[28]. Jahid, M. K. A. S. R. (2024a). Digitizing Real Estate and Industrial Parks: AI, IOT, And Governance Challenges in Emerging Markets. *International Journal of Business and Economics Insights*, *4*(1), 33-70. https://doi.org/10.63125/kbqs6122

[29]. Jahid, M. K. A. S. R. (2024b). Social Media, Affiliate Marketing And E-Marketing: Empirical Drivers For Consumer Purchasing Decision In Real Estate Sector Of Bangladesh. *American Journal of Interdisciplinary Studies*, *5*(02), 64-87. https://doi.org/10.63125/7c1ghy29

[30]. Jahid, M. K. A. S. R. (2025a). AI-Driven Optimization And Risk Modeling In Strategic Economic Zone Development For Mid-Sized Economies: A Review Approach. *International Journal of Scientific Interdisciplinary Research*, *6*(1), 185-218. https://doi.org/10.63125/31wna449

[31]. Jahid, M. K. A. S. R. (2025b). The Role Of Real Estate In Shaping The National Economy Of The United States. *ASRC Procedia: Global Perspectives in Science and Scholarship*, *1*(01), 654–674. https://doi.org/10.63125/34fgrj75

[32]. Karim, S., Shafiullah, M., & Naeem, M. A. (2024). When one domino falls, others follow: A machine learning analysis of extreme risk spillovers in developed stock markets. *International Review of Financial Analysis*, *93*, 103202. https://doi.org/10.1016/j.irfa.2024.103202

[33]. Khandani, A. E., Kim, A. J., & Lo, A. W. (2010). Consumer credit-risk models via machine-learning algorithms. *Journal of Banking & Finance*, *34*(11), 2767-2787. https://doi.org/10.1016/j.jbankfin.2010.06.001

[34]. Kozodoi, N., Jacob, J., & Lessmann, S. (2021). Fairness in credit scoring: Assessment, implementation and profit implications. *European Journal of Operational Research*, *297*(3), 1083-1094. https://doi.org/10.1016/j.ejor.2021.06.023

[35]. Lessmann, S., Baesens, B., Seow, H.-V., & Thomas, L. C. (2015). Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research. *European Journal of Operational Research*, *247*(1), 124-136. https://doi.org/10.1016/j.ejor.2015.05.030

[36]. Lokanan, M. E. (2023). Predicting mobile money transaction fraud using machine learning algorithms. *Applied AI Letters*, *4*, e85. https://doi.org/10.1002/ail2.85

[37]. Makridakis, S., Spiliotis, E., & Assimakopoulos, V. (2018). The M4 Competition: Results, findings, conclusion and way forward. *International Journal of Forecasting*, *34*(4), 802-808. https://doi.org/10.1016/j.ijforecast.2018.06.001

[38]. Malekipirbazari, M., & Aksakalli, V. (2015). Risk assessment in social lending via random forests. *Expert Systems with Applications*, *42*(10), 4621-4631. https://doi.org/10.1016/j.eswa.2015.02.001

[39]. Md Arifur, R., & Sheratun Noor, J. (2022). A Systematic Literature Review of User-Centric Design In Digital Business Systems: Enhancing Accessibility, Adoption, And Organizational Impact. *Review of Applied Science and Technology*, *1*(04), 01-25. https://doi.org/10.63125/ndjkpm77

[40]. Md Hasan, Z., Sheratun Noor, J., & Md. Zafor, I. (2023). Strategic role of business analysts in digital transformation tools, roles, and enterprise outcomes. *American Journal of Scholarly Research and Innovation*, *2*(02), 246-273. https://doi.org/10.63125/rc45z918

[41]. Md Ismail, H., Md Mahfuj, H., Mohammad Aman Ullah, S., & Shofiul Azam, T. (2025). Implementing Advanced Technologies For Enhanced Construction Site Safety. *American Journal of Advanced Technology and Engineering Solutions*, *1*(02), 01-31. https://doi.org/10.63125/3v8rpr04

[42]. Md Ismail Hossain, M. A. B., amp, & Mousumi Akter, S. (2023). Water Quality Modelling and Assessment Of The Buriganga River Using Qual2k. *Global Mainstream Journal of Innovation, Engineering & Emerging Technology*, *2*(03), 01-11. https://doi.org/10.62304/jieet.v2i03.64

[43]. Md Jakaria, T., Md, A., Zayadul, H., & Emdadul, H. (2025). Advances In High-Efficiency Solar Photovoltaic Materials: A Comprehensive Review Of Perovskite And Tandem Cell Technologies. *American Journal of Advanced Technology and Engineering Solutions*, *1*(01), 201-225. https://doi.org/10.63125/5amnvb37

[44]. Md Nur Hasan, M. (2024). Integration Of Artificial Intelligence And DevOps In Scalable And Agile Product Development: A Systematic Literature Review On Frameworks. *ASRC Procedia: Global Perspectives in Science and Scholarship*, *4*(1), 01–32. https://doi.org/10.63125/exyqj773

[45]. Md Nur Hasan, M. (2025). Role Of AI And Data Science In Data-Driven Decision Making For It Business Intelligence: A Systematic Literature Review. *ASRC Procedia: Global Perspectives in Science and Scholarship*, *1*(01), 564-588. https://doi.org/10.63125/n1xpym21

[46]. Md Nur Hasan, M., Md Musfiqur, R., & Debashish, G. (2022). Strategic Decision-Making in Digital Retail Supply Chains: Harnessing AI-Driven Business Intelligence From Customer Data. *Review of Applied Science and Technology*, *1*(03), 01-31. https://doi.org/10.63125/6a7rpy62

[47]. Md Redwanul, I., & Md. Zafor, I. (2022). Impact of Predictive Data Modeling on Business Decision-Making: A Review Of Studies Across Retail, Finance, And Logistics. *American Journal of Advanced Technology and Engineering Solutions*, *2*(02), 33-62. https://doi.org/10.63125/8hfbkt70

[48]. Md Rezaul, K., & Md Mesbaul, H. (2022). Innovative Textile Recycling and Upcycling Technologies For Circular Fashion: Reducing Landfill Waste And Enhancing Environmental Sustainability. *American Journal of Interdisciplinary Studies*, *3*(03), 01-35. https://doi.org/10.63125/kkmerg16

[49]. Md. Sakib Hasan, H. (2022). Quantitative Risk Assessment of Rail Infrastructure Projects Using Monte Carlo Simulation And Fuzzy Logic. *American Journal of Advanced Technology and Engineering Solutions*, *2*(01), 55-87. https://doi.org/10.63125/h24n6z92

[50]. Md. Tarek, H. (2022). Graph Neural Network Models For Detecting Fraudulent Insurance Claims In Healthcare Systems. *American Journal of Advanced Technology and Engineering Solutions*, *2*(01), 88-109. https://doi.org/10.63125/r5vsmv21

[51]. Md. Zafor, I. (2025). A Meta-Analysis Of AI-Driven Business Analytics: Enhancing Strategic Decision-Making In SMEs. *Review of Applied Science and Technology*, *4*(02), 33-58. https://doi.org/10.63125/wk9fqv56

[52]. Md.Kamrul, K., & Md Omar, F. (2022). Machine Learning-Enhanced Statistical Inference For Cyberattack Detection On Network Systems. *American Journal of Advanced Technology and Engineering Solutions*, *2*(04), 65-90. https://doi.org/10.63125/sw7jzx60

[53]. Md.Kamrul, K., & Md. Tarek, H. (2022). A Poisson Regression Approach to Modeling Traffic Accident Frequency in Urban Areas. *American Journal of Interdisciplinary Studies*, *3*(04), 117-156. https://doi.org/10.63125/wqh7pd07

[54]. Medeiros, M. C., Vasconcelos, G. F. R., Veiga, Á., & Zilberman, E. (2021). Forecasting inflation in a data-rich environment: The benefits of machine learning methods. *Journal of Business & Economic Statistics*, *39*(1), 98-119. https://doi.org/10.1080/07350015.2019.1637745

[55]. Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., & Gebru, T. (2019). *Model cards for model reporting* Proceedings of the Conference on Fairness, Accountability, and Transparency,

[56]. Moin Uddin, M. (2025). Impact Of Lean Six Sigma On Manufacturing Efficiency Using A Digital Twin-Based Performance Evaluation Framework. *ASRC Procedia: Global Perspectives in Science and Scholarship*, *1*(01), 343-375. https://doi.org/10.63125/z70nhf26

[57]. Moin Uddin, M., & Rezwanul Ashraf, R. (2023). Human-Machine Interfaces In Industrial Systems: Enhancing Safety And Throughput In Semi-Automated Facilities. *American Journal of Interdisciplinary Studies*, *4*(01), 01-26. https://doi.org/10.63125/s2qa0125

[58]. Momena, A., & Md Nur Hasan, M. (2023). Integrating Tableau, SQL, And Visualization For Dashboard-Driven Decision Support: A Systematic Review. *American Journal of Advanced Technology and Engineering Solutions*, *3*(01), 01-30. https://doi.org/10.63125/4aa43m68

[59]. Mubashir, I., & Abdul, R. (2022). Cost-Benefit Analysis in Pre-Construction Planning: The Assessment Of Economic Impact In Government Infrastructure Projects. *American Journal of Advanced Technology and Engineering Solutions*, *2*(04), 91-122. https://doi.org/10.63125/kjwd5e33

[60]. Mubashir, I., & Jahid, M. K. A. S. R. (2023). Role Of Digital Twins and Bim In U.S. Highway Infrastructure Enhancing Economic Efficiency And Safety Outcomes Through Intelligent Asset Management. *American Journal of Advanced Technology and Engineering Solutions*, *3*(03), 54-81. https://doi.org/10.63125/hftt1g82

[61]. Omar Muhammad, F., & Md.Kamrul, K. (2022). Blockchain-Enabled BI For HR And Payroll Systems: Securing Sensitive Workforce Data. *American Journal of Scholarly Research and Innovation*, *1*(02), 30-58. https://doi.org/10.63125/et4bhy15

[62]. Reduanul, H., & Mohammad Shoeb, A. (2022). Advancing AI in Marketing Through Cross Border Integration Ethical Considerations And Policy Implications. *American Journal of Scholarly Research and Innovation*, *1*(01), 351-379. https://doi.org/10.63125/d1xg3784

[63]. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). *"Why should I trust you?": Explaining the predictions of any classifier* Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining,

[64]. Rudin, C. (2019). Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, *1*(5), 206-215. https://doi.org/10.1038/s42256-019-0048-x

[65]. Sadhwani, A., Giesecke, K., & Sirignano, J. (2021). Deep learning for mortgage risk. *Journal of Financial Econometrics*, *19*(2), 313-368. https://doi.org/10.1093/jjfinec/nbaa025

[66]. Sanjai, V., Sanath Kumar, C., Maniruzzaman, B., & Farhana Zaman, R. (2023). Integrating Artificial Intelligence in Strategic Business Decision-Making: A Systematic Review Of Predictive Models. *International Journal of Scientific Interdisciplinary Research*, *4*(1), 01-26. https://doi.org/10.63125/s5skge53

[67]. Sanjai, V., Sanath Kumar, C., Sadia, Z., & Rony, S. (2025). AI And Quantum Computing For Carbon-Neutral Supply Chains: A Systematic Review Of Innovations. *American Journal of Interdisciplinary Studies*, *6*(1), 40-75. https://doi.org/10.63125/nrdx7d32

[68]. Sheratun Noor, J., & Momena, A. (2022). Assessment Of Data-Driven Vendor Performance Evaluation in Retail Supply Chains: Analyzing Metrics, Scorecards, And Contract Management Tools. *American Journal of Interdisciplinary Studies*, *3*(02), 36-61. https://doi.org/10.63125/0s7t1y90

[69]. Sirignano, J., & Cont, R. (2019). Universal features of price formation in financial markets: Perspectives from deep learning. *Quantitative Finance*, *19*(9), 1449-1459. https://doi.org/10.1080/14697688.2019.1622295

[70]. Tahmina Akter, R., Debashish, G., Md Soyeb, R., & Abdullah Al, M. (2023). A Systematic Review of AI-Enhanced Decision Support Tools in Information Systems: Strategic Applications In Service-Oriented Enterprises And Enterprise Planning. *Review of Applied Science and Technology*, *2*(01), 26-52. https://doi.org/10.63125/73djw422

[71]. Yang, G., Liu, X., & Li, B. (2023). Anti-money laundering supervision by intelligent algorithm. *Computers & Security*, *132*, 103344. https://doi.org/10.1016/j.cose.2023.103344

[72]. Yeh, I.-C., & Lien, C.-H. (2009). The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. *Expert Systems with Applications*, *36*(2), 2473-2480. https://doi.org/10.1016/j.eswa.2007.12.020