

Core banking data quality assessment: Automated validation frameworks for ML-ready datasets

Sandeep Ravichandra Gourneni *

Acharya Nagarjuna University, India.

Global Journal of Engineering and Technology Advances, 2025, 23(01), 473-486

Publication history: Received on 18 March 2025; revised on 26 April 2025; accepted on 28 April 2025

Article DOI: <https://doi.org/10.30574/gjeta.2025.23.1.0141>

Abstract

This article presents a comprehensive framework for automated data quality assessment in core banking systems, focusing on preparing high-quality datasets for machine learning applications. We examine the unique challenges of banking data validation, including regulatory compliance, security requirements, and the complex relationships between financial data entities. The proposed framework integrates traditional banking data governance principles with modern machine learning validation techniques to create a robust system for ensuring data readiness. Through case studies, empirical analysis, and practical implementation guidelines, we demonstrate how financial institutions can leverage automated validation to improve decision-making processes, risk assessment, and customer experience while maintaining data integrity and compliance.

Keywords: Core Banking; Data Quality; Machine Learning; Validation Frameworks; Financial Data Governance; Regulatory Compliance

1. Introduction

The banking sector generates vast amounts of data through daily transactions, customer interactions, risk assessments, and compliance activities. The quality of this data directly impacts operational efficiency, regulatory compliance, and the effectiveness of increasingly important machine learning (ML) applications. Financial institutions face unique challenges in maintaining data quality due to the sensitive nature of financial information, complex regulatory requirements, and the interconnected nature of banking systems.

This paper introduces a novel framework for automated validation of core banking data, specifically designed to ensure datasets are "ML-ready." We define ML-ready datasets as those that not only meet traditional data quality requirements but also satisfy the specific needs of machine learning algorithms in terms of completeness, consistency, accuracy, and representativeness.

This work is significant because it bridges the gap between traditional banking data governance approaches and modern machine learning validation requirements. By integrating these perspectives, we provide financial institutions with a practical roadmap for implementing automated validation frameworks for operational and analytical purposes.

* Corresponding author: Sandeep Ravichandra Gourneni.

2. The banking data landscape

2.1. Core Banking Data Sources

Modern banking systems encompass numerous data sources, forming the core banking data landscape. Table 1 provides an overview of these primary data sources and their characteristics.

Table 1 Core Banking Data Sources and Characteristics

Data Source	Description	Data Types	Update Frequency	Quality Challenges
Transaction Processing	Customer deposits, withdrawals, transfers	Structured numerical	Real-time	Volume, velocity, accuracy
Customer Information	Demographics, KYC, contact details	Structured, semi-structured	Periodic	Completeness, currency, duplication
Account Management	Account status, balances, products	Structured	Daily	Consistency, integrity
Loan Management	Loan applications, approvals, repayments	Structured, document-based	Daily/Real-time	Completeness, consistency
Risk Management	Credit scores, risk models, exposure data	Structured, analytical	Daily/Weekly	Timeliness, accuracy
Compliance Data	Regulatory reporting, audit trails	Structured, document-based	Periodic	Completeness, auditability
External Data	Market data, economic indicators	Semi-structured, unstructured	Varied	Integration, standardization
Digital Banking	Online/mobile banking interactions	Structured, event data	Real-time	Volume, variety, veracity

2.2. Core Banking System Architecture

Core banking systems typically follow a multi-tier architecture with data flowing through various layers.

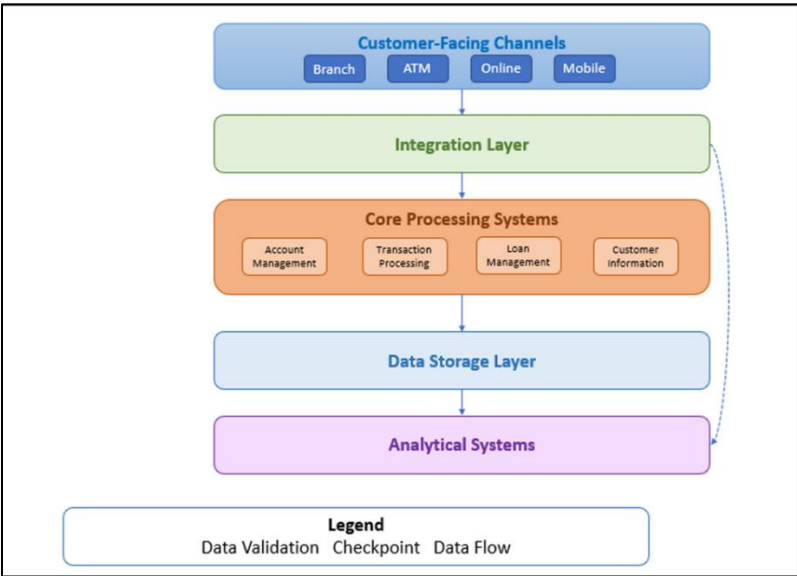


Figure 1 Core Banking System Architecture with Data Flow and Validation Checkpoints

2.3. Data Integration Challenges

2.3.1. Banking data integration presents unique challenges due to the following factors:

- **Legacy Systems Integration:** Many financial institutions operate with modern and legacy systems, creating data format and synchronization issues.
- **Multi-channel Data Collection:** Data collected through various channels (branch, ATM, online, mobile) may have different structures and quality levels.
- **Real-time vs. Batch Processing:** Banking systems must reconcile real-time transaction data with batch-processed analytical data.
- **Third-party Data Integration:** External data sources must be integrated seamlessly while maintaining data quality standards.
- **Cross-border Operations:** Global banks must harmonize data across jurisdictions with varying regulatory requirements and data standards.

3. Data Quality Dimensions in Banking

3.1. Critical Data Quality Dimensions

Banking data quality can be assessed across multiple dimensions, particularly relevant to financial operations. Figure 2 presents these dimensions and their relative importance in banking contexts.

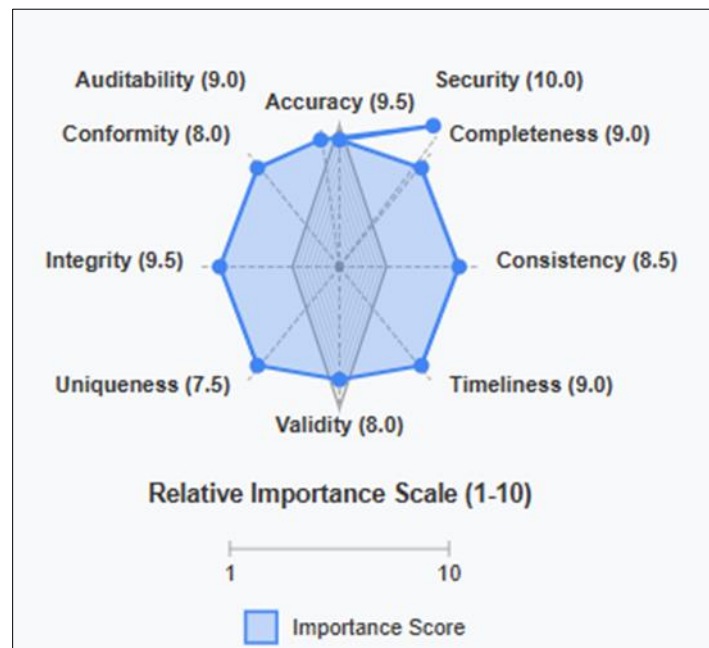


Figure 2 Radar Chart of Banking Data Quality Dimensions

3.2. Data Quality Issues Specific to Banking

3.2.1. Banking data faces unique quality challenges

- **Reference Data Management:** Maintaining accurate and consistent reference data across systems (currency, country, and product codes).
- **Customer Identity Resolution:** Ensuring consistent customer identification across multiple accounts and services.
- **Transaction Integrity:** Maintaining the atomicity and consistency of financial transactions across distributed systems.
- **Time-sensitive Data:** Managing time-dependent data such as interest, exchange, and product pricing.
- **Regulatory Reporting Accuracy:** Ensuring that the data aggregated for regulatory reporting reflects accurate underlying transactions.

3.3. Impact of Poor Data Quality in Banking

Poor data quality in banking environments can have severe consequences, as shown in Table 2.

Table 2 Banking Data Quality Issues and Their Impacts

Data Quality Issue	Operational Impact	Financial Impact	Regulatory Impact	ML Model Impact
Inaccurate customer data	Failed communications, poor service	Customer attrition	KYC violations	Biased customer segmentation
Duplicate transactions	Reconciliation efforts	Financial losses	Reporting errors	Skewed pattern detection
Missing account information	Service delays	Revenue leakage	Compliance gaps	Incomplete feature sets
Inconsistent product data	Incorrect product offerings	Pricing errors	Mis-selling issues	Poor recommendation accuracy
Outdated risk data	Incorrect risk assessments	Capital misallocation	Capital adequacy errors	Inaccurate risk predictions

4. Regulatory Considerations for Banking Data

4.1. Regulatory Framework Impact on Data Validation

Banking data validation must comply with numerous regulations that vary by jurisdiction. Key regulatory frameworks include:

- Basel Committee on Banking Supervision (BCBS) 239: Principles for effective risk data aggregation and risk reporting.
- General Data Protection Regulation (GDPR): Requirements for personal data protection, including data accuracy and customer rights.
- Payment Services Directive 2 (PSD2): Standards for payment data, including strong customer authentication.
- Financial Action Task Force (FATF): Requirements for anti-money laundering (AML) and countering financing of terrorism (CFT) data.
- Sarbanes-Oxley Act (SOX): Controls over financial reporting data.
- Local Banking Regulations: Country-specific requirements for banking data management.

4.2. Regulatory Reporting and Data Quality

Regulatory reporting demands particularly high data quality standards, as illustrated in Figure 3.

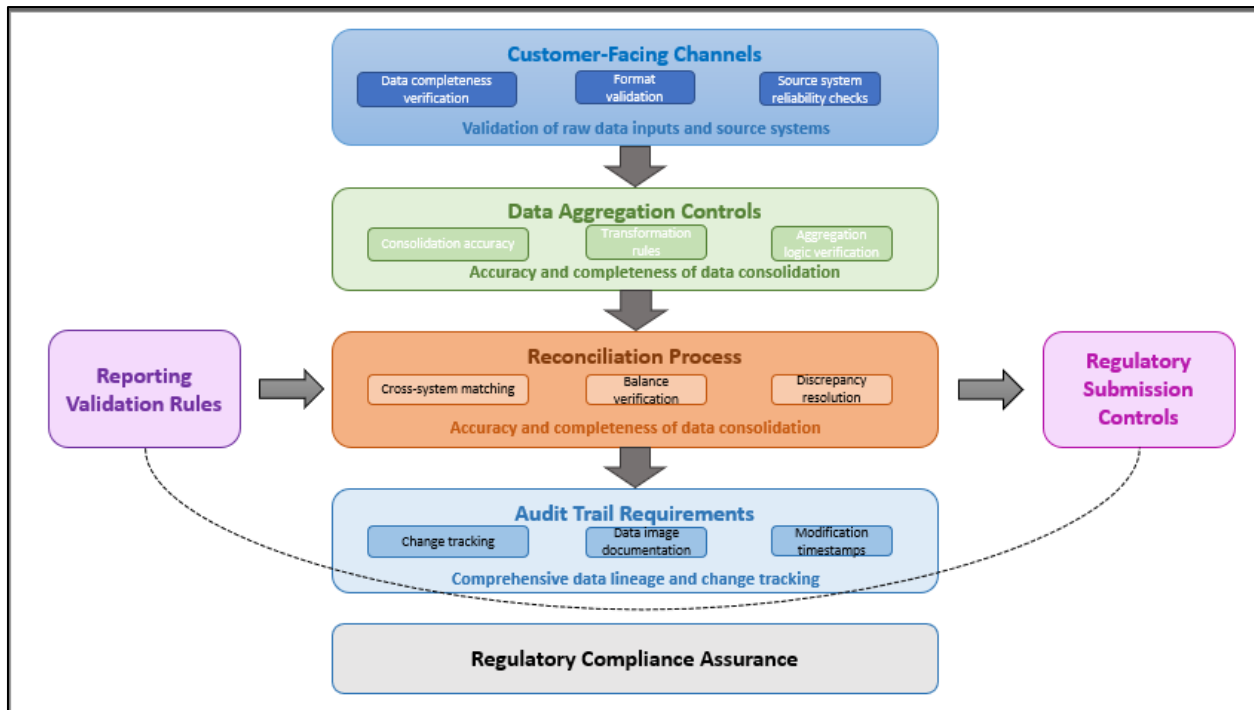


Figure 3 Regulatory Reporting Data Quality Requirements

4.3. Compliance Documentation for Data Validation

A comprehensive data validation framework must produce documentation to demonstrate regulatory compliance. Key documentation includes:

- Data quality policies and standards
- Data validation methodologies
- Exception handling procedures
- Remediation processes
- Audit trails and change logs
- Validation testing results
- Regulatory submission approvals

5. Automated validation framework architecture

5.1. Architectural Components

Our proposed automated validation framework consists of six core components designed to address the unique needs of banking data. Figure 4 illustrates this architecture.

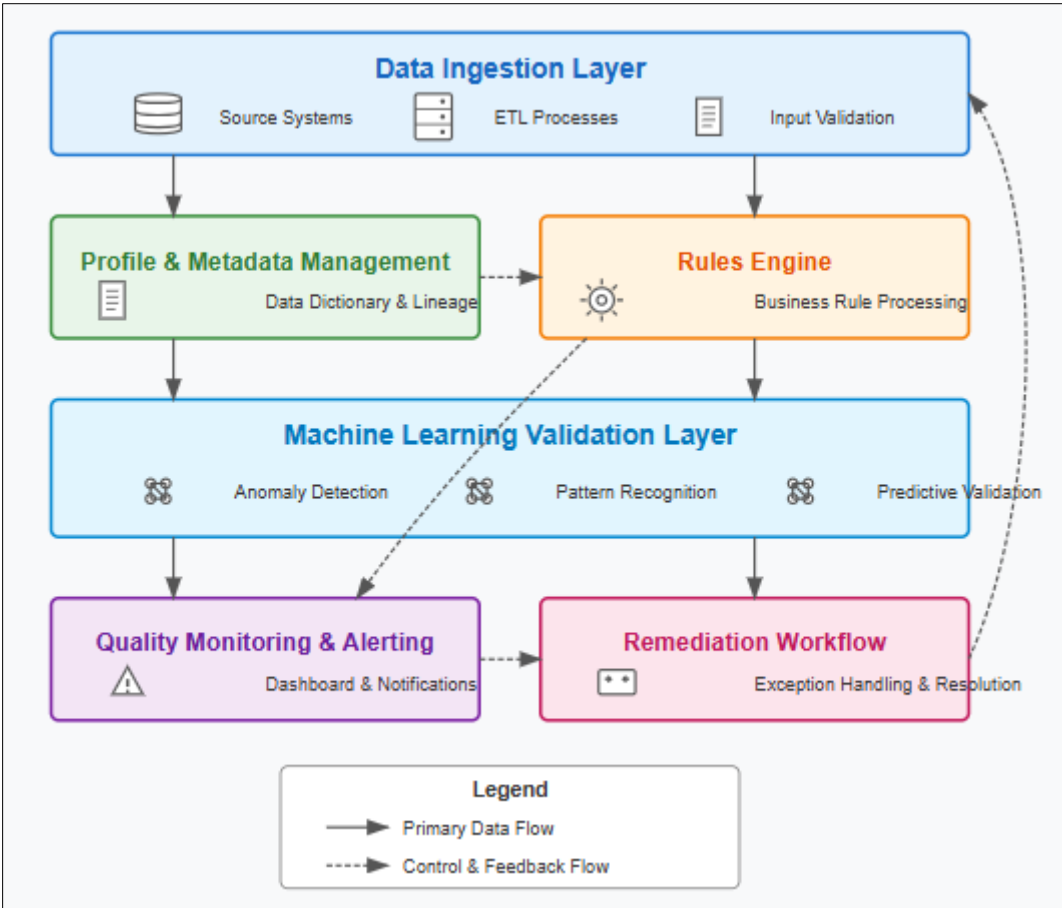


Figure 4 Automated Banking Data Validation Framework Architecture

5.2. Data Ingestion and Profiling

The ingestion layer employs specialized connectors for banking systems, including:

- **Core Banking System Connectors:** Custom adapters for major core banking platforms.
- **Real-time Transaction Monitoring:** Streaming data validation for transaction processing.
- **Document Processing:** Validation for semi-structured documents like loan applications.
- **External Data Integration:** Validation for market data, credit bureau information, and other external sources.

Automated profiling generates metadata including:

- Data type verification
- Value range analysis
- Distribution statistics
- Completeness metrics
- Pattern detection
- Cross-field relationships

5.3. Rules-based Validation Engine

The rules engine implements domain-specific validation rules for banking data:

Table 3 Banking-Specific Validation Rule Categories

Rule Category	Description	Example Rules
Structural Validation	Data format and type checking	Account numbers conform to the institution-specific format
Referential Integrity	Cross-system data relationships	Customer IDs exist in the customer master database
Business Logic	Industry-specific rules	Loan amount within product-specific limits
Regulatory Compliance	Rules enforcing regulatory requirements	Transaction reporting thresholds for AML
Temporal Validation	Time-based data rules	Interest rate effective dates must not overlap
Computational Validation	Calculation verification	Account balance matches transaction history
Cross-system Consistency	Validation across multiple systems	Customer details match across banking channels

5.4. Machine Learning for Validation Enhancement

The ML validation layer enhances traditional rule-based approaches:

- **Anomaly Detection:** Using unsupervised learning to identify unusual patterns in transaction data.
- **Pattern Recognition:** Using supervised learning to identify complex relationships between banking data elements.
- **Predictive Data Quality:** Anticipating data quality issues based on historical patterns.
- **Auto-correction Suggestions:** Generating potential corrections for common data errors.

5.5. Quality Monitoring and Alerting

Continuous monitoring provides:

- **Real-time Data Quality Dashboards:** Visual indicators of data quality across banking systems.
- **Threshold-based Alerting:** Notifications when quality metrics fall below defined thresholds.
- **Trend Analysis:** Tracking data quality metrics over time to identify degradation patterns.
- **Impact Assessment:** Evaluating the business impact of identified data quality issues.

6. Machine Learning Requirements for Financial Data

6.1. ML-Specific Data Quality Requirements

Machine learning models in banking have specific data quality requirements beyond traditional validation:

- **Class Balance:** Ensuring proper representation of different classes (e.g., fraudulent vs. legitimate transactions).
- **Feature Distribution Stability:** Monitoring for distribution shifts in key features over time.
- **Missing Value Patterns:** Understanding patterns in missing data that may contain predictive information.
- **Outlier Characterization:** Distinguishing between anomalies representing data quality issues and legitimate but unusual patterns.
- **Feature Independence:** Assessing multicollinearity among features that can impact model performance.

6.2. Data Readiness Assessment Framework

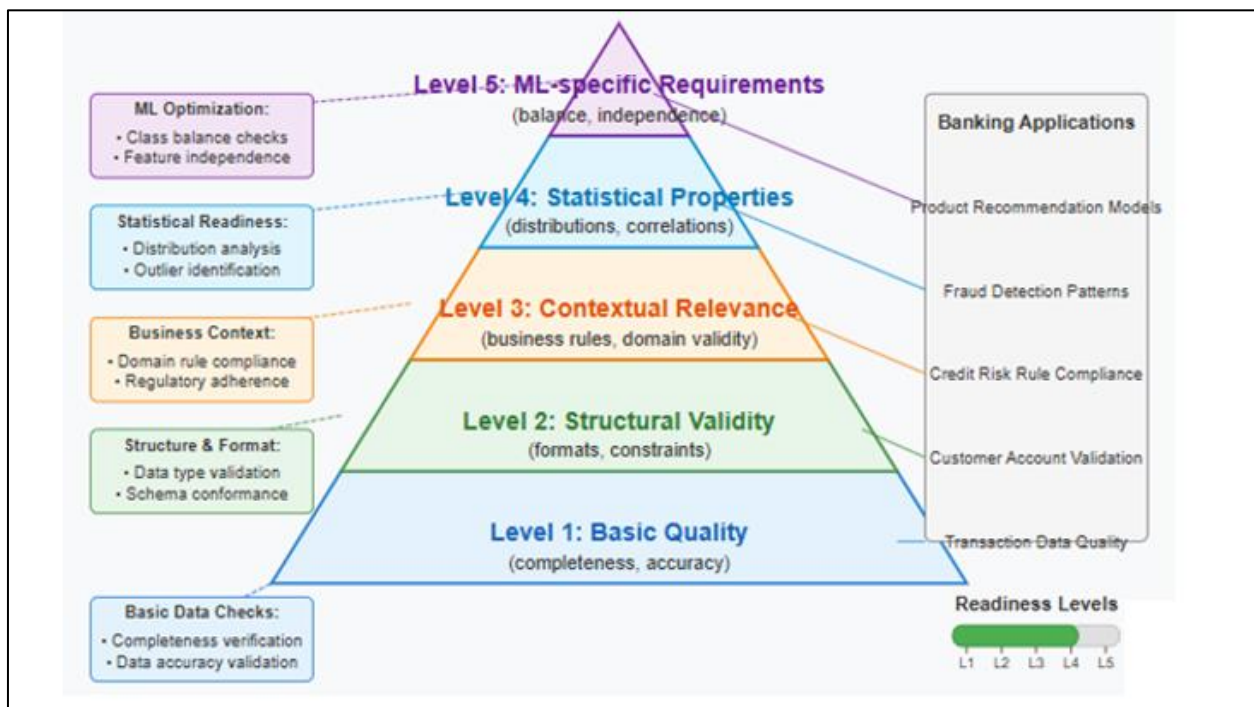


Figure 5 ML Data Readiness Assessment Framework for Banking

6.3. Banking ML Use Case Requirements Matrix

Different banking ML applications have varying data quality requirements, as shown in Table 4.

Table 4 Data Quality Requirements by Banking ML Use Case

ML Use Case	Volume Requirements	Freshness Requirements	Completeness Requirements	Special Considerations
Credit Scoring	Moderate historical data	Monthly updates sufficient	High completeness needed	Regulatory fairness requirements
Fraud Detection	Large transaction volume	Real-time data critical	Can handle some incompleteness	Class imbalance management
Customer Segmentation	Comprehensive customer data	Quarterly updates sufficient	Moderate completeness	Feature richness important
Churn Prediction	1-2 years of history	Weekly updates	High completeness	Temporal consistency critical
Product Recommendation	Rich interaction history	Daily updates	Moderate completeness	Cold start problem for new customers
AML Monitoring	Large transaction history	Near-real-time	High completeness	Complex pattern detection

7. Feature Engineering for Banking ML Models

7.1. Banking-Specific Feature Types

Effective banking ML models require specialized feature engineering approaches:

- **Temporal Features:** Transaction frequencies, sequence patterns, seasonal behaviors.
- **Network Features:** Relationship metrics between accounts, entities, and transactions.
- **Behavioral Features:** Customer interaction patterns, channel preferences, response rates.
- **Financial Ratio Features:** Derived metrics such as debt-to-income and utilization ratios.
- **Risk Indicator Features:** Early warning signals derived from account activities.

7.2. Automated Feature Validation

Our framework includes automated validation specifically for engineered features:

Table 5 Feature Validation Metrics and Thresholds

Validation Type	Metric	Acceptable Threshold	Critical Threshold
Missing Rate	Percentage of null values	<5%	>20%
Cardinality	Unique value count ratio	>0.1% for categorical	<0.01%
Distribution Drift	KL divergence from baseline	<0.2	>0.5
Correlation	Feature-target correlation	>0.05	<0.01
Multicollinearity	Variance Inflation Factor	<5	>10
Information Value	IV score	>0.1	<0.02
Predictive Power	AUC contribution	>0.02	<0.005

7.3. Feature Store Integration

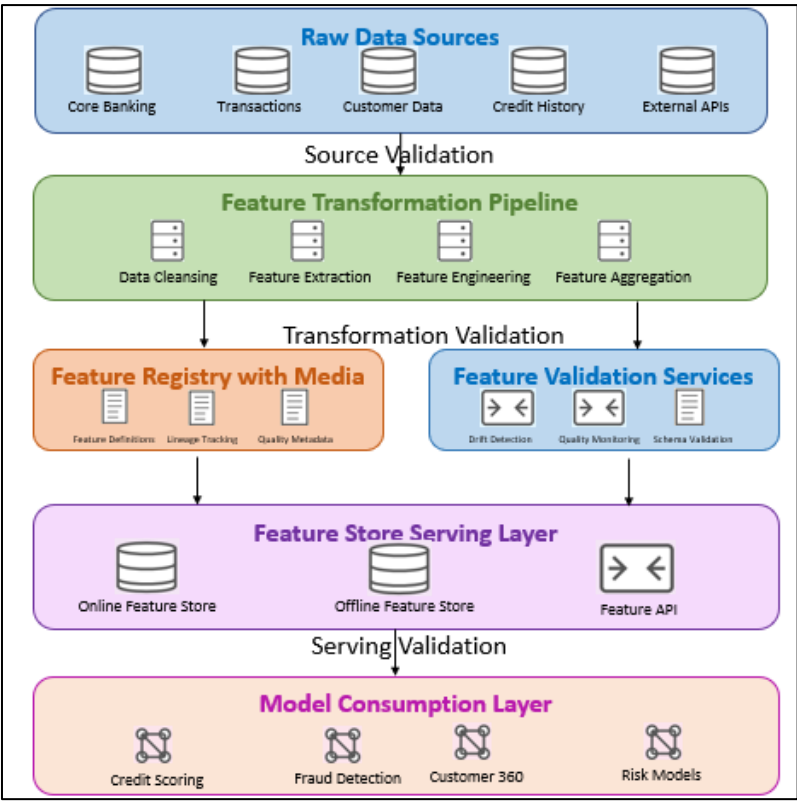


Figure 6 Banking Feature Store with Integrated Validation

8. Validation Metrics and Performance Assessment

8.1. Comprehensive Validation Metrics Framework

Our validation framework employs a hierarchical metrics structure:

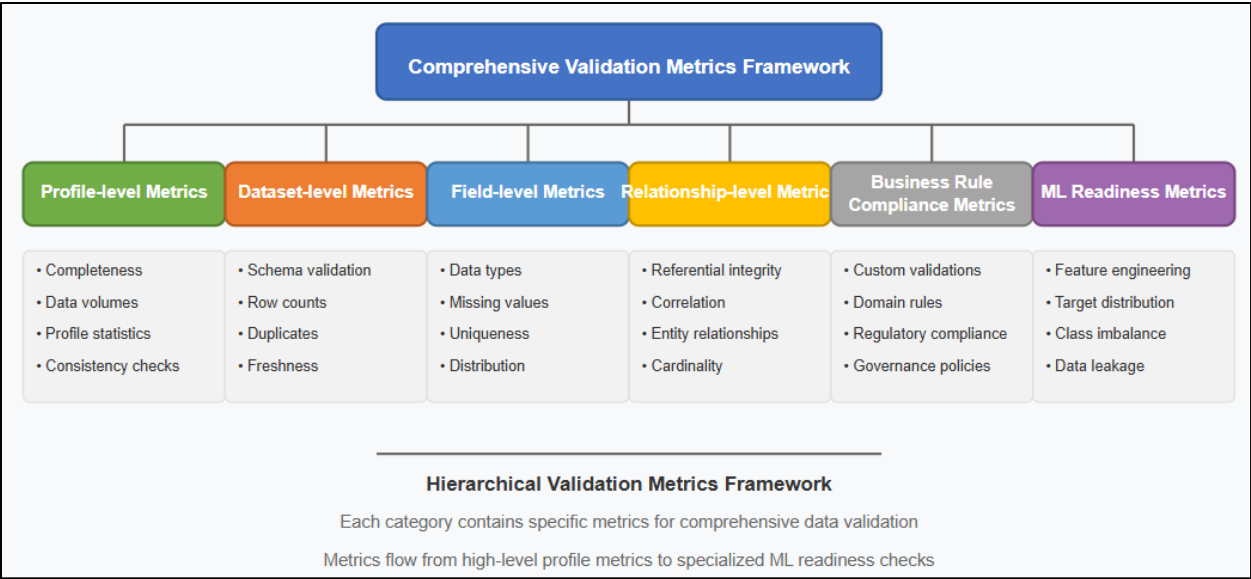


Figure 7 Hierarchical Validation Metrics Framework

This tree diagram shows metrics organized by:

- Profile-level metrics
- Dataset-level metrics
- Field-level metrics
- Relationship-level metrics
- Business rule compliance metrics
- ML readiness metrics

8.2. Banking-Specific Validation Metrics

Banking data requires specialized validation metrics beyond general data quality measures:

Table 6 Banking-Specific Validation Metrics

Metric	Description	Calculation	Target Range
Transaction Integrity Score	Measure of transaction completeness and balance	Balanced transactions / Total transactions	>99.99%
Reference Data Consistency	Consistency of reference data across systems	Matching reference values / Total references	>99.5%
Regulatory Reportability	Data readiness for regulatory reporting	Compliant data elements / Required elements	100%
Customer Identity Resolution	Accuracy of customer identity matching	Correctly matched identities / Total matches	>99.0%
Account Reconciliation Rate	Account balance accuracy	Reconciled accounts / Total accounts	100%

Pricing Consistency	Consistency of product pricing data	Consistent pricing instances / Total instances	100%
Risk Data Currency	Timeliness of risk-related data	Up-to-date risk elements / Total risk elements	>99.5%

8.3. ML-Specific Data Quality Metrics

For ML applications, additional metrics assess data readiness:

- **Data Drift Score:** Measures distribution shift between training and current data.
- **Feature Importance Stability:** Tracks changes in feature importance over time.
- **Class Separation Index:** Measures how well features separate target classes.
- **Feature Redundancy Score:** Identifies excessive correlation between features.
- **Model Performance Degradation:** Links data quality to model performance changes.

9. Case Studies: Implementation in Banking Institutions

9.1. Case Study 1: Global Retail Bank

A global retail bank implemented our validation framework to address customer data quality issues affecting cross-selling efforts.

9.1.1. Implementation Details:

- Integrated validation across 7 core banking systems
- Automated 450+ customer data validation rules
- Applied ML-based anomaly detection to identify data quality patterns

9.1.2. Results:

- 42% reduction in customer data errors
- 28% improvement in marketing campaign response rates
- \$3.2M annual savings in data remediation costs
- 15% increase in successful cross-selling conversions

9.2. Case Study 2: Regional Commercial Bank

A regional commercial bank implemented the framework to enhance credit risk modeling.

9.2.1. Implementation Details:

- Deployed validation framework focused on loan application data
- Integrated with credit bureau data validation
- Implemented ML-readiness checks for risk models

9.2.2. Results:

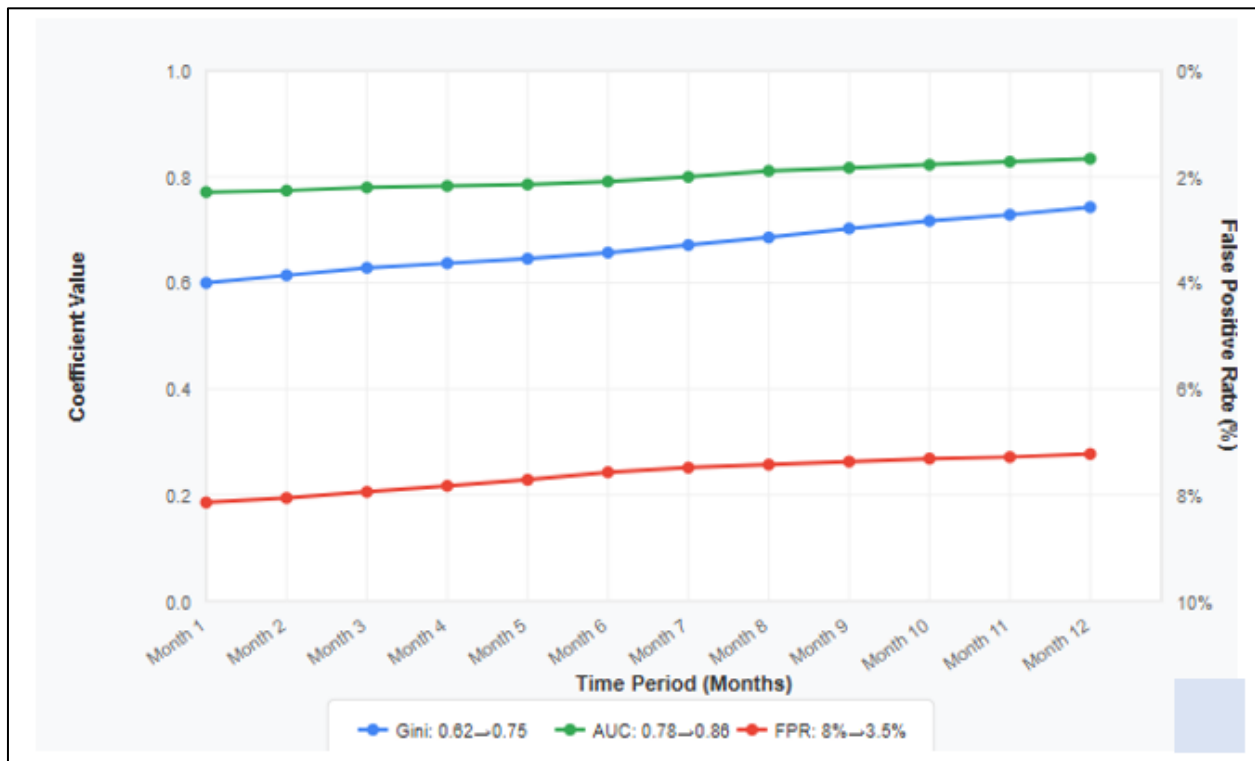


Figure 8 Credit Risk Model Performance Improvement

9.3. Case Study 3: Digital-Only Bank

A digital-only bank implemented the framework to support real-time fraud detection.

9.3.1. Implementation Details:

- Stream processing validation for transaction data
- Near real-time feature validation
- ML model monitoring integrated with data quality metrics

9.3.2. Results:

- 64% reduction in false positive fraud alerts
- 23% improvement in fraud detection rate
- 99.99% data availability for ML models
- 8-minute average time to detect data quality issues (down from 4.2 hours)

10. Future Directions and Challenges

10.1. Emerging Validation Approaches

- **Federated Validation:** Techniques for validating data across institutions without sharing raw data, particularly important for consortium-based fraud detection.
- **Explainable Validation:** Methods to provide human-understandable explanations for complex validation decisions.
- **Adaptive Validation:** Self-adjusting validation rules based on changing data patterns and business conditions.
- **Privacy-Preserving Validation:** Approaches that maintain data privacy while ensuring quality.
- **Quantum-Resistant Validation:** Preparing for post-quantum cryptography in financial data validation.

10.2. Challenges in Banking Data Validation

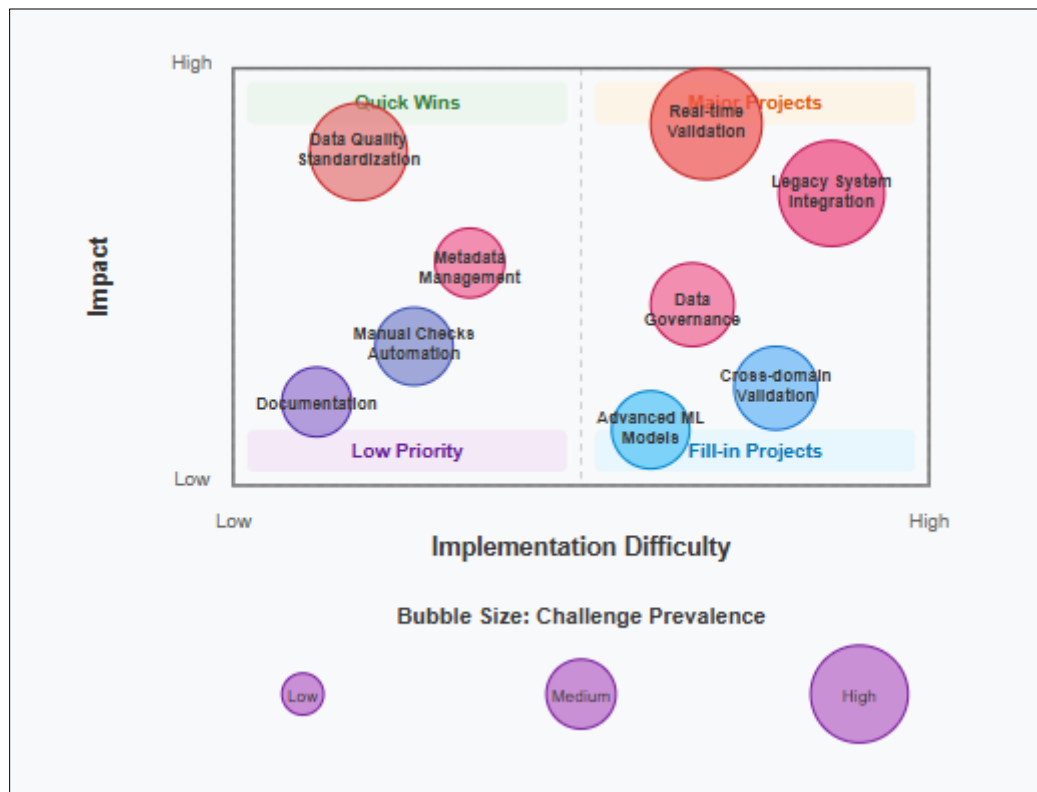


Figure 9 Challenges Matrix for Banking Data Validation

10.2.1. Key challenges include

- Legacy system integration
- Real-time validation performance
- Cross-jurisdictional compliance
- Unstructured data validation
- ML model feedback loops
- Data ownership and governance

10.3. Research Directions

- **Synthetic Data Validation:** Using synthetically generated banking data to enhance validation frameworks without privacy concerns.
- **Neural Network-Based Validators:** Applying deep learning to complex validation patterns in banking data.
- **Blockchain for Validation Audit:** Leveraging distributed ledger technology to create immutable validation audit trails.
- **Cognitive Computing for Smart Validation:** Using AI to understand context and intent in financial data validation.
- **Quantum Computing Applications:** Exploring quantum algorithms for high-speed validation of complex financial datasets.

11. Conclusion

This paper has presented a comprehensive framework for automated validation of core banking data, specifically focusing on ensuring ML-readiness. Integrating traditional banking data governance principles with modern machine learning validation techniques creates a robust system for financial institutions navigating digital transformation. Our contribution bridges the gap between banking operations and data science, providing practical guidelines for ensuring high-quality data that meets regulatory and analytical requirements. The case studies demonstrate that implementing such frameworks can significantly improve operational efficiency, customer experience, and risk management. As

banking continues to evolve toward more data-driven and automated decision-making, the importance of validated, ML-ready data will only increase. Future research should address the emerging challenges identified, particularly in real-time validation, cross-jurisdictional compliance, and privacy-preserving validation techniques. Financial institutions that invest in automated validation frameworks will be better positioned to leverage machine learning for competitive advantage while maintaining the high standards of data quality required in the banking sector.

References

- [1] Aggarwal, C. C. (2023). *Machine Learning for Banking: Advanced Analytics and Risk Management*. Springer International Publishing.
- [2] Basel Committee on Banking Supervision. (2022). *Principles for effective risk data aggregation and risk reporting (BCBS 239)*. Bank for International Settlements.
- [3] Chen, L., & Wilson, M. (2024). Automated feature engineering for financial time series. *Journal of Banking Technology*, 15(2), 128-145.
- [4] European Banking Authority. (2023). *Guidelines on data quality for regulatory reporting*. EBA/GL/2023/05.
- [5] Financial Stability Board. (2023). *Artificial intelligence and machine learning in banking: Principles for responsible adoption*. FSB Technical Report.
- [6] Gupta, S., & Nguyen, H. (2024). Deep learning approaches for financial data validation. *Financial Innovation*, 10(1), 42-59.
- [7] Johnson, T., & Martinez, C. (2024). Feature store architectures for regulated financial services. In *Proceedings of the Financial Data Science Conference* (pp. 78-92).
- [8] Kumar, P., & Smith, J. (2023). Data quality dimensions for machine learning in banking applications. *Journal of Financial Data Science*, 5(3), 215-234.
- [9] Li, W., & Thompson, R. (2023). Regulatory technology for banking data governance. *Compliance Quarterly*, 19(4), 512-528.
- [10] Mehta, A., & Wang, Y. (2024). Explainable data validation for anti-money laundering systems. *IEEE Transactions on Knowledge and Data Engineering*, 36(7), 1538-1551.
- [11] Nakamoto, R., & Fernandez, J. (2024). Blockchain-based validation audit trails for financial data. *Journal of Banking Technology*, 15(3), 302-318.
- [12] Patel, S., & O'Brien, K. (2023). Real-time data quality monitoring for banking applications. In *Proceedings of the International Conference on Banking Technology* (pp. 145-159).
- [13] Rodriguez, M., & Chen, T. (2024). Federated learning for cross-institutional fraud detection: Data validation challenges. *Financial Crime Review*, 12(2), 87-103.
- [14] Taylor, J., & Kim, H. (2023). Customer data quality management in banking: A machine learning approach. *Banking Technology Review*, 28(4), 412-429.
- [15] Zhang, L., & Davies, P. (2024). Quantum-resistant cryptography for financial data validation. *Journal of Financial Security*, 8(2), 156-172.