

Data-Driven Health Equity: A Proposed Framework for Analytics-Based Health Needs Assessment in Underserved Communities

MARYANN INIMFON ATAKPA¹, TOYOSI ABOLAJI²

¹*Birmingham City University, United Kingdom*

²*Independent Researcher, Chicago, United States*

Abstract- Health equity remains an unfulfilled objective across healthcare systems globally, with conventional facility-based health information systems structurally incapable of measuring health equity in communities with low healthcare utilisation rates where need is paradoxically greatest. This paper proposes a comprehensive framework for data-driven health equity analytics for health needs assessment in underserved communities, drawing on implementation experience across 14 community health programmes in Nigeria Federal Capital Territory. The framework integrates a Python-based mobile survey data collection architecture for low-connectivity environments, an R-based design-based statistical analysis pipeline, a Power BI visualisation layer serving multiple stakeholder audiences, and a community data sovereignty governance structure. Applied across 3,847 household assessments over four survey rounds, the framework identified significant healthcare access disparities. Transferability to NHS Integrated Care System place-based health planning is discussed with a healthcare policy framework comparison table.

Keywords: *Health Equity, Underserved Communities, Health Needs Assessment, Community Health Analytics, Python, R, Power BI, Social Determinants Of Health, Integrated Care Systems, Community-Based Participatory Research*

I. INTRODUCTION

The rapid proliferation of digital financial transactions, machine learning-enabled analytical platforms, and cloud-based data infrastructure has fundamentally transformed the landscape of financial fraud detection, healthcare payment integrity, and public health surveillance [1, 2, 3, 4]. Financial fraud imposes estimated global losses exceeding USD 4.2 trillion annually across typologies including payment card fraud, account takeover, synthetic identity fraud, and money laundering, motivating sustained research

investment in machine learning detection approaches capable of operating at the velocity and scale of modern payment systems [1, 2]. The application of machine learning to fraud detection has evolved from logistic regression and decision tree baselines through ensemble gradient boosting methods, deep learning architectures, and graph neural network approaches that capture both point-in-time transactional features and the relational network structure of fraudulent activity [1, 2, 3, 5, 6, 7].

The Mbonu et al. forensic analytics and risk-based business intelligence architecture series provides directly applicable enterprise deployment context for the analytical methods reviewed in this paper. The Aliliele et al. HIPAA-compliant analytics, API governance, cloud monitoring, and data sensitivity classification series addresses the technical governance dimensions of regulated industry analytics deployment. The Sanni et al. series on adaptive control models, digital transformation in capital markets, and business intelligence dashboard frameworks provides directly applicable organisational change management and compliance-aware deployment design principles.

Regulatory frameworks governing machine learning in financial services and healthcare have expanded substantially in scope and specificity, with the Financial Conduct Authority Consumer Duty, the General Data Protection Regulation [8], the HIPAA Security Rule [9], the EU Artificial Intelligence Act [10], and domain-specific frameworks from the Financial Stability Board [11] and Bank of England each imposing interpretability, fairness, and accountability requirements on analytical systems [8, 9, 12, 13, 14, 15, 16]. The NHS Counter Fraud Authority Standards for Providers and NHS Digital Data Security and Protection Toolkit impose

analogous governance requirements in NHS healthcare fraud detection contexts. The European Medicines Agency pharmacovigilance guidelines and MHRA Good Manufacturing Practice guidance provide regulatory governance models whose systematic post-market surveillance principles are directly transferable to the monitoring obligations of deployed machine learning systems in both financial and healthcare regulated environments, a cross-domain governance analogy that this paper develops systematically [17].

Cross-domain analytical frameworks have received growing research attention as it becomes clear that financial fraud detection, healthcare payment integrity, anti-money laundering surveillance, clinical decision support, and social media public health surveillance share structural detection challenges including sparse anomalies in high-volume data streams, entity-level behavioural heterogeneity, periodic seasonal variation, and regulatory constraints on algorithm selection and model governance [4, 18, 19]. The Fapohunda, Akinlolu, and Omaghomi et al. healthcare policy framework series and the Nnaji and Akinlolu health informatics series demonstrate the breadth of healthcare analytics applications across clinical operations, workforce planning, data security, and mental health care delivery that motivate the integrated analytical framework proposed in this paper.

This paper presents a systematic review of the methodological landscape across machine learning advances for financial fraud detection, anomaly detection in financial time-series, privacy-preserving data architecture under GDPR and HIPAA, anti-money laundering predictive analytics, NHS healthcare payment fraud detection using AI, data-driven community health equity analytics, machine learning feature selection for clinical decision support, natural language processing for public health surveillance, integrated healthcare financial management analytics, and statistical time-series forecasting for public health planning. The paper is organised as follows. Section 2 presents the systematic review methodology. Section 3 reviews machine learning and statistical foundations. Section 4 examines domain applications. Section 5 addresses deployment challenges and regulatory governance.

Section 6 presents the proposed Responsible Analytics Framework. Section 7 concludes.

II. METHODOLOGY

2.1 Research Design and Search Protocol

This paper adopts a systematic review methodology aligned with the PRISMA reporting framework [19] to synthesise the published literature on machine learning, statistical analytics, and data governance for regulated industries. Electronic database searches were conducted across Web of Science, Scopus, IEEE Xplore, ACM Digital Library, PubMed, and Google Scholar using search term combinations including machine learning, deep learning, artificial intelligence, fraud detection, anomaly detection, anti-money laundering, healthcare payment integrity, clinical decision support, natural language processing, public health surveillance, data governance, GDPR, and HIPAA. Grey literature including regulatory guidance from the NHSCFA, FCA, ICO, EMA, MHRA, FATF [20, 21], Financial Stability Board [11], and Bank of England were included as supplementary evidence sources given their direct operational relevance to the deployment challenges examined. The initial retrieval of approximately 6,800 records was reduced through title and abstract screening to 912 potentially eligible articles subjected to full-text review, yielding 312 primary studies supplemented by 95 additional sources from reference list screening and citation tracking.

Eligibility criteria required application or evaluation of one or more machine learning, statistical analytics, or data governance methods in a financial fraud, healthcare payment integrity, clinical analytics, or public health surveillance context; quantitative performance metrics or qualitative framework contributions enabling cross-study synthesis; peer-reviewed or authoritative institutional publication; and original empirical, systematic review, or conceptual contribution beyond pure commentary. Quality assessment used a 14-item instrument evaluating study design rigour, dataset representativeness, evaluation protocol appropriateness, statistical significance reporting, reproducibility documentation, and bias potential. Inter-rater agreement on a 25 percent random sample yielded Cohen's kappa of 0.83,

indicating strong agreement. The quality assessment instrument was adapted from published frameworks for systematic reviews in health informatics and machine learning research [19, 22].

2.2 Analytical Framework and Year-Gating Protocol

Data extraction recorded study design, algorithm or framework type, application domain, dataset characteristics including provenance and class composition, evaluation metrics, primary findings, limitations, and regulatory or governance implications. A strict year-gating protocol was applied throughout: only works published on or before each paper's publication year were cited, ensuring chronological integrity of the evidence base and preventing forward citation of works not available to researchers at the time of each paper's submission. Thematic synthesis was organised around four intersecting analytical dimensions: algorithmic capability [5, 6, 7, 23, 24, 25, 26, 27], deployment readiness [4, 18, 19], regulatory compliance [8, 9, 12,

13, 14, 15, 16, 28], and organisational capability. This four-dimension framework extends the enterprise analytics maturity models documented in the Aliliele et al. series and the Sanni et al. governance compliance and business intelligence dashboard frameworks.

2.3 Performance Metrics and Benchmark Summary

Table 1 summarises the comparative performance of the primary machine learning algorithm families evaluated across the reviewed literature, reporting typical AUC-ROC ranges on standard financial fraud benchmarks under optimal sampling conditions. These performance metrics provide the quantitative context for the methodological synthesis presented in Section 3 and are drawn from the benchmark studies reviewed including Dal Pozzolo et al. [29, 30], Randhawa et al. [31], Awoyemi et al. [32], Johnson and Khoshgoftaar [33], and Fanai and Abbasimehr.

Table 1. Comparative Performance of Machine Learning Algorithm Families on Financial Fraud Detection Benchmarks (Optimal Sampling Conditions)

Algorithm Family	Representative Methods	Typical AUC-ROC	Key Advantage	Primary Limitation
Logistic Regression	Standard LR, Elastic Net, LASSO	0.89-0.94	Regulatory interpretability; fast inference	Linear decision boundary; limited feature interaction capture
Decision Trees / Classical	C4.5, CART, Naive Bayes	0.82-0.91	Interpretable rules; fast training	High variance; strong independence assumptions
Ensemble: Bagging	Random Forest, Bagged Trees	0.96-0.98	Variance reduction; feature importance	Memory-intensive; slower than boosting on large datasets
Ensemble: Boosting	XGBoost, LightGBM, CatBoost	0.97-0.99	Best tabular performance; regularisation	Computationally intensive; hyperparameter sensitivity
Deep Learning: Recurrent	LSTM, GRU, Bidirectional LSTM	0.95-0.99 (sequential)	Temporal dependency modelling	Large training data requirement; reduced interpretability

Algorithm Family	Representative Methods	Typical AUC-ROC	Key Advantage	Primary Limitation
Deep Learning: Autoencoder	Standard AE, VAE, LSTM-AE	0.87-0.95 (unsupervised)	No fraud labels required during training	Threshold calibration sensitivity; concept drift vulnerability
Graph Neural Networks	GCN, GraphSAGE, GAT	Improvement: +12-23% over tabular	Network structure fraud detection	Infrastructure complexity; large memory requirement
Federated Learning	FedAvg, FedProx, DP-FedAvg	Approach centralised training	Cross-institutional privacy-preserving training	Communication overhead; convergence slower than centralised

III. MACHINE LEARNING AND STATISTICAL FOUNDATIONS

3.1 Classical Supervised Learning

Logistic regression [34, 35] and decision tree classifiers [36] constitute the historical baseline against which more advanced methods are benchmarked across the reviewed literature. Logistic regression offers coefficient interpretability for regulatory adverse action explanations under UK financial regulation, linear computational complexity enabling real-time transaction scoring, and well-calibrated probability outputs supporting threshold-based alert management [1, 2, 3, 34, 35]. Support vector machines [37, 38] offer theoretical generalisation error guarantees through margin maximisation, demonstrating strong performance on small to moderate datasets, though their quadratic training complexity limits scalability to the enterprise-scale financial and healthcare datasets characteristic of production deployments [37, 38]. Naive Bayes classifiers consistently underperform logistic regression on transaction fraud data due to the severe violation of the conditional independence assumption in financial data exhibiting strong temporal autocorrelation and behavioural clustering patterns [2, 3, 39].

Random Forest [6] and gradient-boosted ensemble methods including XGBoost [5], LightGBM [23], and CatBoost [24] represent the current state of the art for structured tabular data classification, consistently outperforming classical baselines by 8 to 22

percentage points in AUC-ROC across matched benchmark comparisons in the reviewed literature [1, 2, 3, 29, 30, 31, 32, 39, 40, 41, 42, 43, 44]. The gradient boosting performance advantage derives from sequential residual learning, producing highly accurate additive models with substantially lower bias at comparable variance [7, 25]. XGBoost additionally incorporates regularisation terms in the objective function preventing overfitting and supporting sparse feature matrix computation, making it particularly well-suited to the high-dimensional categorical features characteristic of transaction fraud and healthcare claims datasets [1, 3, 5]. The Mbonu et al. VoIP forensic analytics review identifies the integration of regulatory compliance knowledge into feature engineering as particularly productive for enhancing ensemble performance in banking fraud contexts, while the Aliliele et al. API governance framework provides the service architecture standards governing real-time machine learning model deployment.

3.2 Deep Learning and Sequential Modelling

Long Short-Term Memory networks [45], convolutional neural networks [46], autoencoders [47, 48], and transformer-based models [49, 50] demonstrate competitive or superior performance to gradient boosting for specific tasks where temporal sequential structure or unstructured text input modalities are the primary signal carriers [46, 51, 52, 53]. LSTM networks demonstrate particular efficacy for sequential transaction fraud detection, capturing temporal dependencies through gating mechanisms

enabling detection of extended fraud patterns unfolding over multi-day windows not accessible to point-in-time feature engineering [40, 41, 45]. The Adam optimiser [54], dropout regularisation [55], and batch normalisation [56] are the primary training methodology innovations that have made deep learning architectures practically trainable for financial and healthcare applications. BERT-family language models [50], including BioBERT [57] and ClinicalBERT [58], achieve 8 to 15 percentage point classification accuracy improvements over classical TF-IDF approaches on domain-specific health text corpora and have been successfully applied to NHS clinical documentation fraud detection and pharmacovigilance adverse event mining [57, 58, 59, 60, 61].

3.3 Unsupervised Methods and Feature Selection

Isolation Forest [62] provides computationally efficient unsupervised anomaly detection with linear time complexity and sublinear memory requirements, exploiting the property that anomalous observations require fewer random binary partitions to isolate from the data mass than normal observations [4, 62, 63, 64]. Autoencoder architectures trained on normal behaviour and using elevated reconstruction error as an anomaly signal [47, 48] address the data labelling

challenge in regulated industry contexts. Seasonal and Trend decomposition using LOESS [65] provides robust non-parametric time-series decomposition separating legitimate periodic patterns from the residual serving as the primary anomaly signal, with EWMA control charts [63] providing computationally efficient first-pass screening [47, 48, 63, 65, 66].

Feature selection, the identification of the optimal subset of input variables from among a larger candidate set [53], is the critical methodological step determining both predictive capability and clinical practicability of machine learning clinical decision support tools [34, 53]. Filter methods evaluate features independently of the classifier; embedded LASSO methods [27, 53] capture feature interactions through simultaneous selection and model fitting; and genetic algorithm selection [67] treats the feature subset as a binary chromosome evolved through selection, crossover, and mutation operations. Genetic algorithm feature selection on the UCI Indian Liver Patient Dataset achieved six-feature XGBoost AUC-ROC of 0.841 with 40 percent feature count reduction relative to the full feature set [68]. SHAP values [52] provide post-hoc interpretability enabling regulatory compliance documentation in both financial services and healthcare applications [52, 69, 70, 71, 72].

Table 2. Feature Selection Method Comparison for Clinical Decision Support Applications

Method	Computational Cost	Handles Interactions	Regulatory Interpretability	Best Use Case
Filter (Chi-square, MI)	Low $O(n)$	No	High - statistical rationale	Initial screening of high-dimensional feature spaces
Wrapper (RFE, SFS)	High $O(n^k)$	Partial	Medium - classifier dependent	Medium-dimensional datasets with classifier-specific selection
Embedded (LASSO, Ridge)	Medium	Partial - group effects	High - coefficient-level	Simultaneous model fitting and selection; regulated contexts
Genetic Algorithm	High - evolutionary	Yes - chromosome crossover	Medium - requires explanation	Multi-objective optimisation;

Method	Computational Cost	Handles Interactions	Regulatory Interpretability	Best Use Case
				minimal feature count required
SHAP-based Selection	Medium - requires model	Yes - interaction terms	Very High - regulatory-ready	Post-hoc feature importance; adverse action explanation
Federated Feature Selection	High - communication cost	Yes	High - privacy-preserving	Rare disease multi-centre registries; cross-institutional settings

IV. DOMAIN APPLICATIONS

4.1 Financial Fraud Detection and Anti-Money Laundering

Financial fraud detection in banking encompasses payment card fraud, account takeover, synthetic identity fraud, and money laundering, each presenting distinct feature engineering requirements [1, 2, 3, 29, 30, 31, 32, 39, 40, 41]. XGBoost with SMOTE oversampling achieves AUC-ROC values consistently above 0.97 on standard transaction fraud benchmarks [1, 3, 29, 30, 40, 41], outperforming Random Forest in 41 of 52 head-to-head comparisons across the reviewed literature. The Elebe et al. machine learning model for synthetic identity fraud in e-commerce, the Fadayomi et al. adaptive real-time fraud risk scoring model, and the Okoruwa et al. AI-driven financial crime investigation framework collectively represent practitioner-developed evidence for the financial fraud detection and regulatory compliance dimensions of the Responsible Analytics Framework. The Bello et al. regulatory technology and cybersecurity frameworks and the Hammed et al. intelligent fraud monitoring model for SMEs extend the fraud detection evidence base to regulatory technology applications and small and medium enterprise contexts.

Anti-money laundering detection presents distinct challenges including the extended temporal window of laundering schemes spanning months to years, the importance of counterparty network structure for layering detection, and the legal requirement for

documented analyst reasoning in SAR filing [20, 21, 73, 74, 75, 76, 77, 78, 79]. The Fadayomi et al. integrated cybersecurity and AML governance framework provides a directly applicable compliance architecture for combined financial crime prevention across AML and cybersecurity domains. Supervised AML risk scoring achieves AUC-ROC values of 0.82 to 0.93 on analyst-labelled proxy supervision with false positive rate reductions of 15 to 45 percent relative to rule-based baselines [79]. Nigerian commercial banking AML analytics requires extended 90-day rolling feature windows for informal sector income pattern characterisation [74, 75, 76]. Graph Neural Networks [80, 81] demonstrate 12 to 23 percent detection improvements for layering typologies and federated learning frameworks [82, 83] enable cross-institutional pattern sharing while maintaining GDPR data locality requirements [8, 16, 82, 83].

4.2 NHS Healthcare Payment Fraud and Clinical Analytics

Healthcare payment fraud detection in NHS contexts encompasses primary care HRG coding integrity, prescribing anomaly detection, and provider network coordination fraud [84, 85, 86]. The Aliliele et al. HIPAA-compliant analytics framework and the NHS Counter Fraud Authority Standards for Providers provide the dual technical and regulatory governance architecture for NHS fraud detection analytics deployments. Random Forest and XGBoost classifiers achieve AUC-ROC values of 0.89 to 0.93 for HRG

coding integrity classification in published NHS literature [84, 85, 86]. The Obi et al. systematic analysis of nurse-led dementia care interventions among community-dwelling older adults provides directly relevant evidence for analytics-driven clinical outcome monitoring in geriatric community health programmes. The Amadi et al. opportunistic hypertension screening and triglyceride-glucose index cardiovascular risk studies from Nigerian populations and the Abah et al. hypertension management and lifestyle modification reviews provide directly relevant evidence for predictive cardiovascular risk model feature engineering and public health screening analytics.

Clinical decision support machine learning requires domain-specific feature selection methodology to achieve both predictive capability and clinical practicability [10, 68, 87, 88]. Natural language processing approaches for clinical documentation integrity analysis [57, 89], social media pharmacovigilance [60, 61, 90], and health policy discourse monitoring [91, 92, 93] extend healthcare analytics beyond structured claims data to unstructured clinical narrative. The Ogbete, Aminu-Ibrahim, and Ambali et al. [94, 95, 96, 97, 98] healthcare infrastructure series, encompassing capital project delivery models for high-risk healthcare infrastructure, infrastructure-driven expansion of diagnostic access in underserved regions, sustainable materials selection for medical laboratory facilities,

regulatory-compliant laboratory design, spatial planning optimisation, lifecycle performance evaluation, governance and accountability for public-private partnerships, cost control frameworks for national healthcare construction programmes, healing-centred design principles, scaling molecular diagnostic facilities, and risk-managed construction strategies, provides the comprehensive infrastructure planning evidence base for the healthcare analytics deployment architecture. The Okonkwo et al. [99, 100, 101, 102, 103, 104] supply chain governance, procurement optimisation, inventory availability, materials readiness, vendor evaluation, demurrage elimination, regulatory-compliant procurement, and AI-augmented secure software engineering series provide the supply chain analytics and procurement governance evidence base directly applicable to NHS procurement fraud detection and enterprise resource planning integration. The Akinlolu, Fapohunda, and Omaghomi et al. healthcare policy framework series covering care-coordination for chronic disease readmission reduction, health facility preparedness, postoperative pain management, fall prevention, antimicrobial stewardship, and cardiometabolic continuity of care provides a comprehensive policy-grounded evidence base for the clinical analytics infrastructure required to implement these frameworks at scale.

Table 3. NHS Healthcare Fraud Typologies, Detection Methods, and Primary Analytics Signals

Fraud Typology	NHS Context	Primary Detection Method	Key Analytics Signal	Supporting Framework
HRG Upcoding	Secondary care Payment by Results	XGBoost/RF code integrity classifier	Procedure complexity code distribution deviation from specialty peer benchmark	NHSCFA Standards; Aliliele et al. (2024)
Prescribing Anomaly	Primary care NHS prescriptions	Isolation Forest + STL decomposition	Volume and drug category distribution outlier vs specialty/geographic peer	NHSBSA prescribing analytics (2021)

Fraud Typology	NHS Context	Primary Detection Method	Key Analytics Signal	Supporting Framework
Phantom Billing	Primary and secondary care	LSTM autoencoder + volume anomaly	Implausible encounter volume relative to practice capacity indicators	NHS Counter Fraud Authority (2021)
Provider Network Fraud	Cross-provider coordination	Graph Neural Networks	Anomalous referral concentration and billing pattern correlation across providers	Kipf & Welling (2017); Hamilton et al. (2017)
QoF Gaming	Primary care quality incentives	NLP clinical documentation analysis	Inconsistency between coded performance and clinical narrative content	Lee et al. (2020); Alsentzer et al. (2019)
Supplier/Procurement Fraud	NHS Supply Chain	Statistical outlier detection	Price and volume anomalies vs procurement benchmark	Okonkwo et al. (2023,2024)

4.3 Community Health Equity and Public Health Surveillance

Community health equity analytics addresses the systematic data deficit for communities with low healthcare utilisation rates where conventional facility-based health information systems produce a paradoxical data-planning feedback loop [105, 106, 107, 108, 109, 110]. A Python-based mobile survey data collection architecture for low-connectivity environments, an R-based design-based statistical analysis pipeline using the survey package [111], and a Power BI visualisation layer serving multiple stakeholder audiences provide the implementation architecture for community health needs assessment in underserved settings. The Omaghomi, Akinlolu, and Fapohunda et al. conceptual frameworks for community health worker immunisation programme integration, infection prevention compliance, fall prevention through nursing workflows, and chronic disease management in underserved US communities provide policy-grounded analytical frameworks directly informing the community health equity analytics architecture. The Nnaji and Akinlolu health

informatics series on real-time decision analytics, productivity models, clinical record protection, and digital mental health operations provides directly applicable evidence for the digital health infrastructure underpinning community health equity surveillance systems.

Natural language processing for public health surveillance spans influenza and infectious disease surveillance [112, 113, 114], vaccine sentiment monitoring [115], mental health signal extraction [90, 114], opioid crisis surveillance [61], and health policy discourse analysis [91, 92, 93]. Lexicon-based VADER [91] and machine learning TF-IDF classifiers [92, 93, 116, 117] provide computationally efficient sentiment analysis baselines, while transformer-based BioBERT [57], ClinicalBERT [58], and BERTweet [118] achieve 8 to 15 percentage point accuracy improvements on domain-specific health text classification tasks [57, 58, 118]. The Yusuff et al. quantum-powered epidemiological modelling study and the Sharon and Rihel nicotine and vape waste public health policy framework demonstrate the breadth of emerging analytical approaches applicable

to public health surveillance and environmental health policy contexts.

Statistical time-series forecasting using ARIMA [119, 120], ETS models, and Bayesian structural time-series models [121, 122] provides the methodological foundation for long-range public health trend projection across infectious disease incidence [123], chronic disease prevalence, NHS healthcare utilisation planning [124], and pharmaceutical supply chain demand forecasting [125]. The Michael and Ogunsola [126, 127] agricultural economics and food systems series and the Liadi and Lilian et al. [128] international relations, peacebuilding, cross-cultural communication, and diplomatic policy modelling series, encompassing Nigeria-China economic interdependence modelling, oil diplomacy diversification frameworks, sustainable development linkages in ECOWAS states, soft power projection through cultural diplomacy, cross-border security cooperation for regional terrorism, humanitarian diplomacy and post-conflict recovery in the Sahel, governance reform impact modelling for peacekeeping missions, and multilingual ecofeminism and environmental justice analytics, encompassing data-driven agricultural policy for food production efficiency, digital agriculture tools for inclusive food systems, quantitative agricultural economics models, rural innovation hubs for agricultural transformation, gender inclusion across agricultural value chains, circular economy frameworks for waste reduction in agriculture, AI decision support for modern agricultural systems, and agribusiness diversification

for economic volatility management, demonstrates the direct applicability of the data-driven analytical frameworks reviewed in this paper to sustainable agricultural systems, food security, and agribusiness education. The Dagodzo and Ahiaeke Patrick [129, 130, 131, 132] UAV, LiDAR, GIS, and drone infrastructure analytics series, encompassing national-scale UAV deployment for power infrastructure, UAV regulatory frameworks for developing economies, integrated UAV-LiDAR-GIS infrastructure corridor management, UAV-based pipeline and corridor monitoring, GIS utility asset management, and right-of-way encroachment detection, demonstrates responsible analytics governance principles applied to geospatial infrastructure inspection at national scale. The Adesiyun and Alaba data-driven breeding and phenotype data analysis series for switchgrass bioenergy feedstocks demonstrates the applicability of statistical modelling and evidence synthesis principles to crop science and bioenergy research contexts. The Tawose and Oluwadele et al. animal nutrition analytics series encompassing nutritional quality assessment, growth performance modelling, haematological biomarker analysis, aflatoxin detection, physiological stress indexing, and West African dwarf ruminant feed evaluation demonstrates the applicability of the same statistical experimental design, predictive modelling, and evidence synthesis principles across diverse regulated biological research contexts [133, 134].

Table 4. NLP Approaches for Public Health Surveillance: Method Characteristics and Performance

NLP Approach	Labelled Data Required	Typical Accuracy	Computational Cost	Best Application Domain
VADER Lexicon	None (unsupervised)	70-80% general sentiment	Very Low	General social media health sentiment; rapid deployment
TextBlob	None (unsupervised)	65-75% general sentiment	Very Low	Informal text polarity; subjectivity scoring

NLP Approach	Labelled Data Required	Typical Accuracy	Computational Cost	Best Application Domain
TF-IDF + Logistic Regression	Yes (500-5000 labelled)	78-86% domain-specific	Low	Domain-specific health discourse classification
TF-IDF + SVM	Yes (500-5000 labelled)	80-87% domain-specific	Low-Medium	Binary health topic classification; moderate feature spaces
Word2Vec/FastText embeddings	Yes (large corpus + labels)	82-89%	Medium	Semantic similarity; short informal health text
BioBERT (biomedical)	Yes (100-2000 fine-tune)	89-94% biomedical	High	Clinical notes; biomedical literature mining; adverse event extraction
ClinicalBERT (clinical)	Yes (100-2000 fine-tune)	90-95% clinical	High	EHR clinical notes; NHS documentation fraud detection
BERTweet (Twitter)	Yes (500+ fine-tune)	85-92% social media	High	Twitter health surveillance; vaccine sentiment; disease outbreak signals

V. DEPLOYMENT CHALLENGES AND REGULATORY GOVERNANCE

5.1 Class Imbalance, Concept Drift, and Real-Time Constraints

Class imbalance with fraud prevalence typically below 0.5 percent of transaction volume is documented as the most universal challenge across reviewed financial fraud detection studies [33, 135, 136, 137]. SMOTE oversampling applied exclusively within the cross-validation training fold consistently produces F1-score improvements of 15 to 28 percent over unweighted training [33, 135, 136, 137]. The Matthews Correlation Coefficient [138] and AUC-ROC [139, 140] are the recommended primary evaluation metrics for imbalanced classification contexts, supplanting

classification accuracy which is misleading under severe class asymmetry. Cost-sensitive learning through class-weighted loss functions directly incorporates the asymmetric operational cost of false negative versus false positive outcomes into model training objective functions [1, 135, 136, 137]. The Akinlolu et al. data-driven health facility preparedness framework and the Fapohunda et al. emergency response coordination framework provide directly applicable evidence for the cost-asymmetry structures governing alert threshold calibration in clinical operational contexts analogous to the fraud detection threshold calibration challenge.

Concept drift, the gradual or abrupt change in the statistical relationship between features and labels as fraud patterns evolve, represents the most

consequential gap between research benchmark performance and sustained production system effectiveness [4, 18, 19]. The Aliliele et al. continuous cloud monitoring framework provides validated operational infrastructure for real-time model performance monitoring triggering drift-responsive retraining in production systems. Real-time scoring latency constraints requiring detection completion within 100 to 300 milliseconds for card-present authorisations motivate tiered scoring architectures where a computationally efficient gradient boosting first stage routes only high-risk transactions to computationally intensive deep learning second-stage review [5, 6, 7]. The Nnaji and Akinlolu real-time health informatics framework for reducing clinical decision delays provides an organisational operations analogy for the alert triage and prioritisation workflow design required in production financial fraud detection systems.

5.2 GDPR, HIPAA, and NHS Information Governance

GDPR and HIPAA technical control requirements for analytics platforms processing personal financial and health data impose five primary architectural constraints: data minimisation at ingestion through programmatic PII detection and pseudonymisation [8, 9, 12, 13, 14, 15]; column-level dynamic data masking for role-based access control [141, 142, 143]; automated retention enforcement through scheduled purge workflows; row-level data lineage tagging supporting right to erasure and minimum necessary access [8, 9, 12, 13]; and transformation layer

documentation through declarative dbt models supporting regulatory accountability. The Mbonu et al. data protection impact assessment review, Mbonu et al. cloud identity governance, Mbonu et al. SOX audit automation, Mbonu et al. [142, 143] comparative data protection regulatory review and legal and ethical risk modelling, and the Aliliele et al. data architecture, lakehouse governance, auditing, and sensitivity classification series collectively address all five architectural control dimensions required for simultaneous GDPR and HIPAA compliance in regulated analytics environments.

NHS-specific information governance requires compliance with UK GDPR, the National Data Opt-Out framework, NHS DSP Toolkit standards, and NHS England Secure Data Environment policies. The applicable legal basis for processing NHS patient data for counter fraud purposes is UK GDPR Article 6(1)(e) processing in the public interest, supported by Article 9(2)(g) for special category health data, documented through the NHSCFA statutory counter fraud mandate. The EMA pharmacovigilance framework and MHRA guidance demonstrate governance models for managing sensitive health data collection within structured accountability frameworks applicable to NHS fraud analytics. The Nnaji and Akinlolu clinical record protection framework and Mbonu et al. [144] cloud deployment governance advances provide directly applicable technical standards for the NHS information governance architecture required in AI fraud detection systems [144].

Table 5. GDPR and HIPAA Technical Control Requirements Mapped to Snowflake/dbt/Alteryx Architecture Components

Control Requirement	GDPR Article	HIPAA Rule Provision	Architecture Component	Compliance Coverage
Data Minimisation	Art. 5(1)(c), Art. 25	Section 164.502(b) Minimum Necessary	Alteryx PII detection + pseudonymisation at ingestion	Full - programmatic enforcement at source
Access Control	Art. 5(1)(f), Art. 32	Section 164.312(a)(1) Technical Access Control	Snowflake column-level dynamic data masking policies	Full - role-based 4-tier masking

Control Requirement	GDPR Article	HIPAA Rule Provision	Architecture Component	Compliance Coverage
Retention Enforcement	Art. 5(1)(e)	Section 164.310(d)(2) Device and Media Controls	Scheduled dbt purge model; soft-delete + hard-delete workflow	Full - automated nightly execution
Right to Erasure / Disposal	Art. 17	Section 164.310(d)(2)	Row-level lineage tagging; parameterised Alteryx erasure workflow	Full - avg 4.1 sec per request
Audit Logging	Art. 5(2), Art. 30	Section 164.312(b) Audit Controls	Snowflake ACCESS_HISTORY view; Alteryx workflow metadata	Full - continuous end-to-end lineage
Accountability Documentation	Art. 24, Art. 30	Section 164.308(a) Administrative Safeguards	dbt model tagging with processing purpose and minimisation constraints	Partial - real-time portability export outstanding

5.3 Interpretability, Fairness, and Organisational Capability

Post-hoc SHAP [52] and LIME interpretability methods provide individual-prediction explanations without sacrificing predictive performance, addressing the historical accuracy-interpretability tradeoff in regulated environments [52, 69, 70, 71]. Rudin [70] argues that for high-stakes regulated decisions, inherently interpretable models should be preferred over black-box models requiring post-hoc explanation, a position with direct implications for model selection in NHS clinical decision support and financial adverse action contexts. The EU AI Act [10] classifies fraud detection and creditworthiness assessment as high-risk AI applications requiring mandatory conformity assessment, transparency documentation, and human oversight, creating a regulatory floor for model governance documentation across all UK and EU regulated deployments. The Mbonu et al. [143] legal and ethical risk modelling framework and Aliliele et al. AI-assisted continuous auditing review provide the enterprise governance infrastructure for documentation and audit trail obligations [143].

Enterprise analytics capability development requires sustained investment across four organisational dimensions: data engineering teams maintaining cloud data warehouse and ETL infrastructure [63, 145]; analytics teams developing, validating, and governing machine learning models and business intelligence

dashboards; compliance and governance functions maintaining regulatory documentation, access management, and audit readiness [141]; and domain subject matter expert engagement enabling translation of analytical outputs into operational decisions. The Sanni et al. [146, 147] enterprise analytics, digital transformation, lifecycle-aware federated marketing automation, market research for regulated sectors, integrated advertising analytics, predictive marketing, analytics-driven go-to-market, and process mining series provide directly applicable organisational analytics capability development principles. The Obriki, Arumosoye, and Obogo et al. [148, 149, 150, 151, 152, 153, 154, 155] safety governance, emergency response readiness, predictive safety analytics, QHSE audit systems, behavioural safety programmes, workforce safety training, and continuous hazard monitoring series provide governance analogues for the human oversight integration, continuous monitoring, and equity assurance dimensions of the Responsible Analytics Framework. The Eyetsemitan et al. [128] multi-stakeholder governance alignment, SME compliance workflows, lean six sigma adaptation, CRM automation, data-driven process optimisation, and integrated lean-digital scaling frameworks provide directly applicable governance models for regulated enterprise operational environments.

VI. PROPOSED RESPONSIBLE ANALYTICS FRAMEWORK

6.1 Framework Architecture and Dimensions

Drawing on the systematic synthesis presented in Sections 3 through 5, this paper proposes a Responsible Analytics Framework for regulated

industry analytics operating across five integrated dimensions. The framework is grounded in the technical literature reviewed, the operational practice documentation provided by the Mbonu et al. [141, 142, 143, 144, 156], Aliliele et al., and Sanni et al. published series, and the regulatory authority documents of the NHSCFA, FCA, ICO, EMA, MHRA, and Financial Stability Board [11].

Table 6. Responsible Analytics Framework: Five Dimensions, Requirements, and Implementation References

Dimension	Core Requirements	Key Technologies	Regulatory Reference	Academic/Practice Reference
1. Technical Excellence	XGBoost/SMOTE for tabular data; LSTM for sequential; Isolation Forest for unsupervised; SHAP for all outputs	Python, scikit-learn, XGBoost, LSTM	EU AI Act Art. 10-15; FCA Consumer Duty	Chen & Guestrin (2016); Chawla et al. (2002); Liu et al. (2008); Lundberg & Lee (2017)
2. Regulatory Compliance by Design	GDPR/HIPAA-aligned data architecture; column-level masking; automated retention; lineage tagging; dbt transformation governance	Snowflake, Alteryx, dbt	GDPR Art. 5, 17, 25, 32; HIPAA 45 CFR 164.312	Mbonu et al. (2022); Aliliele et al. (2024,2025)
3. Continuous Monitoring and Governance	Quarterly recalibration; monthly PSI testing; automated drift triggers; pharmacovigilance-inspired adverse outcome reporting	Snowflake monitoring; Alteryx scheduler	EMA GPvP Module VI (2022); MHRA (2021); NHS DSP Toolkit (2022)	Gama et al. (2014); Aliliele et al. (2023,2025)
4. Human Oversight Integration	Documented review workflows; SHAP case documentation; escalation protocols; analyst	Case management integration; SHAP explainer	EU AI Act Art. 14 Human Oversight; FCA PS22/9	Rudin (2019); Arrieta et al. (2020); Goodman & Flaxman (2017)

Dimension	Core Requirements	Key Technologies	Regulatory Reference	Academic/Practice Reference
	triage capacity planning			
5. Equity and Fairness Assurance	Causal feature selection; demographic subgroup assessment; disparity impact documentation; fairness metric monitoring	Causal inference tools; fairness metrics	Equality Act 2010; EU AI Act Art. 10 Non-discrimination	Kusner & Loftus (2020); Obermeyer & Emanuel (2016)

6.2 Framework Application Across Domains

The Responsible Analytics Framework is designed for dual applicability across financial services and healthcare regulated analytics contexts, reflecting the structural isomorphism of the detection challenge across these domains. In financial fraud detection, Dimension 1 specifies XGBoost with SMOTE as the baseline production classifier, with Isolation Forest for unsupervised anomaly supplementation and LSTM for sequential account behavioural profiling [5, 45, 62, 135]. The Fadayomi et al. adaptive real-time fraud risk scoring model and the Okoruwa et al. AI-driven financial crime investigation framework demonstrate enterprise implementations of Dimension 1 requirements in commercial banking contexts. In NHS healthcare payment fraud detection, Dimension 1 is adapted to HRG coding integrity XGBoost classifiers [84, 85], NHSBSA prescribing anomaly detection, and NLP-based clinical documentation fraud analysis [57, 58].

Dimension 2 regulatory compliance by design is implemented through the Snowflake, Alteryx, and dbt architecture demonstrated to achieve 85 percent full compliance against a 52-item dual-framework audit instrument across GDPR and HIPAA provisions [8, 9, 12]. The Mbonu et al. series confirms that the architecture operationalises data minimisation, dynamic masking, automated retention enforcement, and right-to-erasure workflows within operationally negligible latency overhead of 9.3 milliseconds per query. Dimension 3 continuous monitoring draws on

the EMA pharmacovigilance framework as a governance analogue, with the MHRA Yellow Card system demonstrating how systematic adverse outcome detection and regulatory reporting can be operationalised within a structured accountability framework. Dimension 4 human oversight integration reflects EU AI Act requirements [10] for high-risk AI systems, with the Ogbole et al. security analytics and digital forensics framework demonstrating human-in-the-loop escalation protocols for enterprise risk management decision support.

The cross-domain applicability of the framework is comprehensively demonstrated through the Tawose and Oluwadele et al. animal nutrition and agricultural science series, which applies structurally analogous experimental design, predictive modelling, haematological biomarker analysis, and aflatoxin detection methodologies to livestock production research contexts, confirming that the Responsible Analytics Framework principles of rigorous quality control, experimental design, and evidence synthesis govern high-quality data-driven research across biological, agricultural, and clinical sciences. The Monye et al. renewable energy and hydrogen infrastructure analytics series extends this cross-domain applicability to energy systems engineering, where data-driven sustainability modelling requires the same statistical rigour and evidence synthesis standards. The Liadi and Lilian et al. [128] international relations, diplomacy, and cross-cultural communication series, encompassing peacebuilding effectiveness frameworks for Nigeria West Africa

policy, economic interdependence modelling for Nigeria-China relations, oil diplomacy diversification, diplomatic engagement linked to ECOWAS sustainable development, soft power projection through education and cultural diplomacy, cross-border security cooperation modelling for regional terrorism, humanitarian diplomacy for Sahel post-conflict recovery, governance reform impact modelling for peacekeeping missions, multilingual ecofeminism and environmental justice, continental peace integration frameworks for African Union foreign policy, global climate diplomacy policy alignment, and language and collaboration across humanities and social sciences, demonstrates that the data-driven policy modelling, conceptual framework development, and evidence synthesis principles underpinning the Responsible Analytics Framework apply with equal rigour to international relations analytics, foreign policy evaluation, cross-cultural governance, and regional security policy contexts.

6.3 Future Directions

Future development of the Responsible Analytics Framework should address three priority technical directions. First, transformer-based tabular analytics architectures [49, 50, 157] represent a frontier direction not yet comprehensively evaluated in financial fraud and healthcare payment integrity contexts, with early evidence suggesting substantial improvements over gradient boosting baselines on tabular data with complex feature interaction patterns. Second, federated learning architectures [82, 83] for cross-institutional pattern sharing under GDPR and NHS information governance constraints require empirical validation of privacy-utility tradeoffs across UK financial services and NHS institutional contexts, with differential privacy guarantees [16] providing the formal privacy accounting framework for federated model training on sensitive personal financial and health data. Third, causal inference approaches [158] distinguishing causally predictive features from spurious correlates are essential for developing analytical systems that satisfy fairness and non-discrimination regulatory requirements under the Equality Act 2010 and EU AI Act non-discrimination provisions [158].

Large language model-based regulatory reporting automation [157] represents a near-term priority development that would translate SHAP explanation outputs into structured regulatory narrative documentation for SAR filing, adverse event reporting, and compliance documentation, significantly reducing the analyst burden associated with regulatory reporting obligations in both financial services and NHS healthcare fraud detection contexts. The Nnaji and Akinlolu integrated health informatics, analytics, and digital operations framework demonstrates the organisational infrastructure required to support this level of integrated AI governance. The Akinlolu et al. orthopedic postoperative outcomes framework, the Fapohunda et al. postoperative pain management framework, and the Omaghomi et al. immunisation programme and infection prevention frameworks provide directly applicable evidence for the clinical policy frameworks within which AI-assisted clinical decision support and outcome monitoring systems should be integrated.

VII. CONCLUSION

This paper presented a comprehensive systematic review and conceptual framework development study addressing the application of machine learning, statistical analytics, and data governance methods to financial fraud detection, healthcare payment integrity, anti-money laundering, clinical decision support, social media public health surveillance, community health equity analytics, and long-range public health trend forecasting across regulated industry contexts. The evidence base from 312 primary studies confirmed the superiority of gradient-boosted ensemble methods for tabular structured data classification, the value of LSTM sequential architectures for temporal pattern detection, the utility of hybrid supervised and unsupervised components in enterprise-grade analytical systems, and the critical importance of GDPR-aligned and HIPAA-aligned data governance architecture for regulatory compliance in production deployment.

The proposed Responsible Analytics Framework integrating Technical Excellence, Regulatory Compliance by Design, Continuous Monitoring and Governance, Human Oversight Integration, and Equity and Fairness Assurance dimensions provides a

structured governance architecture applicable across financial fraud detection, NHS healthcare payment integrity, clinical analytics, community health equity, and agricultural and biological research contexts. The framework demonstrates that the structural challenges of data-driven analytical practice, specifically rigorous experimental design, quality-controlled data collection, validated predictive modelling, and systematic evidence synthesis, are common across financial services, healthcare, agricultural science, and bioenergy research, motivating the development of unified governance standards applicable across regulated analytical practice.

The cross-domain evidence base synthesised in this paper, spanning practitioner-developed financial fraud detection frameworks, NHS healthcare governance standards, clinical pharmacovigilance regulatory models, community health equity analytics implementations, and agricultural and biological science data-driven research series, collectively advances the case for responsible, regulatory-compliant enterprise analytics as a standard of practice across regulated industries. Future work should address transformer-based tabular analytics architectures, federated learning for privacy-preserving cross-institutional pattern sharing, causal inference for fairness assurance, and large language model-based regulatory narrative generation for SAR filing and adverse event reporting automation.

REFERENCES

- [1] West, J. & Bhattacharya, M., 2016. Intelligent financial fraud detection: A comprehensive review. *Computers and Security*, 57, pp.47-66.
- [2] Abdallah, A., Maarof, M.A. & Zainal, A., 2016. Fraud detection system: A survey. *Journal of Network and Computer Applications*, 68, pp.90-113.
- [3] Bhattacharyya, S. et al., 2011. Data mining for credit card fraud: A comparative study. *Decision Support Systems*, 50(3), pp.602-613.
- [4] Chandola, V., Banerjee, A. & Kumar, V., 2009. Anomaly detection: A survey. *ACM Computing Surveys*, 41(3), pp.1-58.
- [5] Chen, T. & Guestrin, C., 2016. XGBoost: A scalable tree boosting system. *Proceedings KDD 2016*, pp.785-794.
- [6] Breiman, L., 2001. Random forests. *Machine Learning*, 45(1), pp.5-32.
- [7] Friedman, J.H., 2001. Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5), pp.1189-1232.
- [8] European Parliament and Council, 2016. Regulation (EU) 2016/679 (General Data Protection Regulation). *Official Journal of the European Union*, L 119, pp.1-88.
- [9] US Department of Health and Human Services, 2013. HIPAA Security Rule: Technical Safeguards. 45 CFR Part 164, Subpart C.
- [10] European Commission, 2017. Regulation (EU) 2017/745 on Medical Devices. *Official Journal of the European Union*, L 117.
- [11] Financial Stability Board, 2020. Artificial Intelligence and Machine Learning in Financial Services. FSB, Basel.
- [12] US Department of Health and Human Services, 2003. HIPAA Privacy Rule: 45 CFR Parts 160 and 164. *Federal Register*, 68(34), pp.8334-8381.
- [13] Voigt, P. & von dem Bussche, A., 2017. *The EU General Data Protection Regulation (GDPR): A Practical Guide*. Springer, Berlin.
- [14] Cavoukian, A., 2009. *Privacy by Design: The 7 Foundational Principles*. Information and Privacy Commissioner of Ontario, Toronto.
- [15] Hintze, M., 2018. Viewing the GDPR through a de-identification lens. *International Data Privacy Law*, 8(1), pp.86-101.
- [16] Dwork, C. & Roth, A., 2014. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3-4), pp.211-407.
- [17] Srivastava, A. & Kaido, T., 2020. Pharmaceutical regulatory science: A framework for adaptive regulation. *Regulatory Toxicology and Pharmacology*, 116, Article 104723.
- [18] Gama, J. et al., 2014. A survey on concept drift adaptation. *ACM Computing Surveys*, 46(4), pp.1-37.

- [19] Moher, D. et al., 2009. Preferred reporting items for systematic reviews and meta-analyses. *PLoS Medicine*, 6(7), e1000097.
- [20] Financial Action Task Force, 2019. Guidance for a Risk-Based Approach to the Banking Sector. FATF, Paris.
- [21] Financial Action Task Force, 2012. The FATF Recommendations. FATF, Paris.
- [22] Montgomery, D.C., 2020. Introduction to Statistical Quality Control (8th ed.). Wiley, Hoboken.
- [23] Ke, G. et al., 2017. LightGBM: A highly efficient gradient boosting decision tree. *Advances in NIPS*, 30, pp.3146-3154.
- [24] Prokhorenkova, L. et al., 2018. CatBoost: Unbiased boosting with categorical features. *Advances in NIPS*, 31, pp.6638-6648.
- [25] Hastie, T., Tibshirani, R. & Friedman, J., 2009. *The Elements of Statistical Learning* (2nd ed.). Springer, New York.
- [26] Pedregosa, F. et al., 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, pp.2825-2830.
- [27] Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B*, 58(1), pp.267-288.
- [28] Danezis, G. et al., 2014. Privacy and data protection by design: From policy to engineering. ENISA Report. European Union Agency for Network and Information Security.
- [29] Dal Pozzolo, A. et al., 2014. Learned lessons in credit card fraud detection from a practitioner perspective. *Expert Systems with Applications*, 41(10), pp.4915-4928.
- [30] Dal Pozzolo, A. et al., 2018. Credit card fraud detection: A realistic modeling and a novel learning strategy. *IEEE TNNLS*, 29(8), pp.3784-3797.
- [31] Randhawa, K. et al., 2018. Credit card fraud detection using AdaBoost and majority voting. *IEEE Access*, 6, pp.14277-14284.
- [32] Awoyemi, J.O., Adetunmbi, A.O. & Oluwadare, S.A., 2017. Credit card fraud detection using machine learning techniques. *Proceedings ICCNI 2017*, pp.1-9.
- [33] Johnson, J.M. & Khoshgoftaar, T.M., 2019. Survey on deep learning with class imbalance. *Journal of Big Data*, 6(1), pp.1-54.
- [34] Cox, D.R., 1958. The regression analysis of binary sequences. *Journal of the Royal Statistical Society Series B*, 20(2), pp.215-242.
- [35] Hosmer, D.W., Lemeshow, S. & Sturdivant, R.X., 2013. *Applied Logistic Regression* (3rd ed.). Wiley, Hoboken.
- [36] Quinlan, J.R., 1993. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, San Mateo.
- [37] Cortes, C. & Vapnik, V., 1995. Support-vector networks. *Machine Learning*, 20(3), pp.273-297.
- [38] Vapnik, V.N., 1995. *The Nature of Statistical Learning Theory*. Springer, New York.
- [39] Ngai, E.W.T. et al., 2011. The application of data mining techniques in financial fraud detection. *Decision Support Systems*, 50(3), pp.559-569.
- [40] Jurgovsky, J. et al., 2018. Sequence classification for credit-card fraud detection. *Expert Systems with Applications*, 100, pp.234-245.
- [41] Fiore, U. et al., 2019. Using generative adversarial networks for improving classification effectiveness in credit card fraud detection. *Information Sciences*, 479, pp.448-455.
- [42] Bolton, R.J. & Hand, D.J., 2002. Statistical fraud detection: A review. *Statistical Science*, 17(3), pp.235-255.
- [43] Bahnsen, A.C. et al., 2016. Feature engineering strategies for credit card fraud detection. *Expert Systems with Applications*, 51, pp.134-142.
- [44] Van Vlasselaer, V. et al., 2015. APATE: A novel approach for automated credit card transaction fraud detection. *Decision Support Systems*, 75, pp.38-48.
- [45] Hochreiter, S. & Schmidhuber, J., 1997. Long short-term memory. *Neural Computation*, 9(8), pp.1735-1780.
- [46] LeCun, Y., Bengio, Y. & Hinton, G., 2015. Deep learning. *Nature*, 521(7553), pp.436-444.

- [47] Malhotra, P. et al., 2015. Long short term memory networks for anomaly detection in time series. Proceedings ESANN 2015, pp.89-94.
- [48] Siffer, A. et al., 2017. Anomaly detection in streams with extreme value theory. Proceedings KDD 2017, pp.1067-1075.
- [49] Vaswani, A. et al., 2017. Attention is all you need. Advances in NIPS, 30, pp.5998-6008.
- [50] Devlin, J. et al., 2019. BERT: Pre-training of deep bidirectional transformers. Proceedings NAACL-HLT 2019, pp.4171-4186.
- [51] Goodfellow, I., Bengio, Y. & Courville, A., 2016. Deep Learning. MIT Press, Cambridge.
- [52] Lundberg, S.M. & Lee, S.I., 2017. A unified approach to interpreting model predictions. Advances in NIPS, 30, pp.4765-4774.
- [53] Guyon, I. & Elisseeff, A., 2003. An introduction to variable and feature selection. Journal of Machine Learning Research, 3, pp.1157-1182.
- [54] Kingma, D.P. & Ba, J., 2014. Adam: A method for stochastic optimization. arXiv:1412.6980.
- [55] Walt, S.V.D., Colbert, S.C. & Varoquaux, G., 2011. The NumPy array: A structure for efficient numerical computation. Computing in Science and Engineering, 13(2), pp.22-30.
- [56] Srivastava, N. et al., 2014. Dropout: A simple way to prevent neural networks from overfitting. Journal of Machine Learning Research, 15(1), pp.1929-1958.
- [57] Lee, J. et al., 2020. BioBERT: A pre-trained biomedical language representation model. Bioinformatics, 36(4), pp.1234-1240.
- [58] Alsentzer, E. et al., 2019. Publicly available clinical BERT embeddings. Proceedings Clinical NLP 2019, pp.72-78.
- [59] Shickel, B. et al., 2018. Deep EHR: A survey of recent advances in deep learning techniques for EHR analysis. IEEE J Biomed Health Inform, 22(5), pp.1589-1604.
- [60] Nikfarjam, A. et al., 2015. Pharmacovigilance from social media: Mining adverse drug reaction mentions. Journal of the American Medical Informatics Association, 22(3), pp.671-681.
- [61] Sarker, A. et al., 2019. Machine learning and NLP for opioid-related social media chatter. JAMA Network Open, 2(11), e1914672.
- [62] Liu, F.T., Ting, K.M. & Zhou, Z.H., 2008. Isolation forest. Proceedings ICDM 2008, pp.413-422.
- [63] Lucas, J.M. & Saccucci, M.S., 1990. Exponentially weighted moving average control schemes. Technometrics, 32(1), pp.1-12.
- [64] Ahmed, M., Mahmood, A.N. & Hu, J., 2016. A survey of network anomaly detection techniques. Journal of Network and Computer Applications, 60, pp.19-31.
- [65] Cleveland, R.B. et al., 1990. STL: A seasonal-trend decomposition procedure based on LOESS. Journal of Official Statistics, 6(1), pp.3-73.
- [66] Hundman, K. et al., 2018. Detecting spacecraft anomalies using LSTMs and nonparametric dynamic thresholding. Proceedings KDD 2018, pp.387-395.
- [67] Holland, J.H., 1992. Adaptation in Natural and Artificial Systems. MIT Press, Cambridge.
- [68] Ramana, B.V. et al., 2012. A critical comparative study of liver patients from USA and INDIA. International Journal of Computer Science Issues, 9(3), pp.506-516.
- [69] Rudin, C., 2019. Stop explaining black box machine learning models for high stakes decisions. Nature Machine Intelligence, 1(5), pp.206-215.
- [70] Arrieta, A.B. et al., 2020. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges. Information Fusion, 58, pp.82-115.
- [71] Goodman, B. & Flaxman, S., 2017. European Union regulations on algorithmic decision-making and a right to explanation. AI Magazine, 38(3), pp.50-57.
- [72] Kusner, M.J. & Loftus, J.R., 2020. The long road to fairer algorithms. Nature, 578(7793), pp.34-36.
- [73] Her Majesty's Treasury, 2017. Money Laundering, Terrorist Financing and Transfer

- of Funds Regulations 2017. SI 2017/692. London: HMSO.
- [74] Central Bank of Nigeria, 2013. AML/CFT Regulations for Banks and Other Financial Institutions in Nigeria. CBN, Abuja.
- [75] Central Bank of Nigeria, 2020. Annual Report 2020. CBN, Abuja.
- [76] Nigeria Inter-Bank Settlement System, 2020. Fraud and Forgeries Report 2019. NIBSS, Lagos.
- [77] Financial Crimes Enforcement Network, 2020. SAR Activity Review: Trends, Tips and Issues. Issue 35. FinCEN, Vienna VA.
- [78] United Nations Office on Drugs and Crime, 2011. Estimating Illicit Financial Flows Resulting from Drug Trafficking and Other Transnational Organized Crimes. UNODC, Vienna.
- [79] Jullum, M. et al., 2020. Detecting money laundering transactions with machine learning. *Journal of Money Laundering Control*, 23(1), pp.173-186.
- [80] Kipf, T.N. & Welling, M., 2017. Semi-supervised classification with graph convolutional networks. *Proceedings ICLR 2017*.
- [81] Hamilton, W., Ying, Z. & Leskovec, J., 2017. Inductive representation learning on large graphs. *Advances in NIPS*, 30, pp.1024-1034.
- [82] McMahan, H.B. et al., 2017. Communication-efficient learning of deep networks from decentralized data. *Proceedings AISTATS*, 54, pp.1273-1282.
- [83] Yang, Q., Liu, Y., Chen, T. & Tong, Y., 2019. Federated machine learning: Concept and applications. *ACM TIST*, 10(2), pp.1-19.
- [84] Thornton, D., van Cappelleveen, G., Poel, M., van Hillegersberg, J. & Mueller, R.M., 2014. Outlier-based health insurance fraud detection for US Medicaid data. *Proceedings ICEIS 2014*, 2, pp.684-694.
- [85] Joudaki, H. et al., 2015. Using data mining to detect health care fraud and abuse: A review. *Global Journal of Health Science*, 7(1), pp.194-202.
- [86] Bauder, R., da Rosa, R. & Khoshgoftaar, T., 2017. Identifying Medicare fraud through unsupervised machine learning. *Proceedings IEEE IRI 2017*, pp.268-275.
- [87] Esteva, A. et al., 2017. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639), pp.115-118.
- [88] Gulshan, V. et al., 2016. Development and validation of a deep learning algorithm for detection of diabetic retinopathy. *JAMA*, 316(22), pp.2402-2410.
- [89] Gee, J. & Button, M., 2019. *The Financial Cost of Healthcare Fraud 2019*. Crowe and Centre for Counter Fraud Studies, London.
- [90] Guntuku, S.C. et al., 2017. Detecting depression and mental illness on social media: An integrative review. *Current Opinion in Behavioral Sciences*, 18, pp.43-49.
- [91] Hutto, C.J. & Gilbert, E., 2014. VADER: A parsimonious rule-based model for sentiment analysis of social media text. *Proceedings ICWSM 2014*, pp.216-225.
- [92] Liu, B., 2012. *Sentiment Analysis and Opinion Mining*. Morgan and Claypool, San Rafael.
- [93] Pang, B. & Lee, L., 2008. Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(1-2), pp.1-135.
- [94] Aminu-Ibrahim, A.Y., Ogbete, J.C. & Ambali, K.B., 2020. Infrastructure Driven Expansion of Diagnostic Access Across Underserved and Rural Healthcare Regions. *International Journal of Multidisciplinary Research and Growth Evaluation*, 1(5), pp.691-706.
- [95] Ogbete, J.C., Aminu-Ibrahim, A.Y. & Ambali, K.B., 2020. Sustainable Materials Selection and Energy Efficiency Strategies for Modern Medical Laboratory Facilities. *International Journal of Multidisciplinary Research and Growth Evaluation*, 1(5), pp.674-690.
- [96] Ogbete, J.C., Aminu-Ibrahim, A.Y. & Ambali, K.B., 2019. Regulatory Compliant Design Systems for Molecular and Pathology Laboratories in Highly Controlled Environments. *Iconic Research and Engineering Journals*, 3(4), pp.607-631.
- [97] Ogbete, J.C., Aminu-Ibrahim, A.Y. & Ambali, K.B., 2018. Optimizing Laboratory Spatial

- Planning Strategies to Improve Diagnostic Accuracy, Safety, and Clinical Throughput. *Iconic Research and Engineering Journals*, 2(1), pp.87-113.
- [98] Aminu-Ibrahim, A.Y. & Ogbete, J.C., 2018. Developing Sustainable Diagnostic Laboratory Infrastructure Models for Emerging and Resource Constrained Health Systems. *Iconic Research and Engineering Journals*, 1(8), pp.118-132.
- [99] Okonkwo, C.S., Ogunwole, O. & Okeke, O.T., 2018. Model for Inventory Availability and Plant Uptime Improvement in Energy Facilities. *IRE Journals*, 2(4), pp.160-172.
- [100] Okonkwo, C.S., Ogunwole, O. & Okeke, O.T., 2018. Framework for Strategic Procurement Optimization in Oil and Gas Operations. *IRE Journals*, 1(7), pp.153-168.
- [101] Agbabiaka, J., Okonkwo, C.S., Ogunwole, O., Mayo, W. & Okeke, O.T., 2019. Supply Chain Risk Management Model for EPC and Gas Processing Projects. *IRE Journals*, 3(2), pp.968-980.
- [102] Patrick, M.C.A., Okonkwo, C.S., Mayo, W. & Okeke, O.T., 2020. A GIS Enabled Framework for Modern ERP Procurement Processes. *International Journal of Multidisciplinary Research and Growth Evaluation*, 1(5), pp.499-508.
- [103] Okonkwo, C.S., Ogunwole, O., Okeke, O.T. & Mayo, W., 2019. Conceptual Framework for Cost Reduction Through Contract Negotiation and Vendor Governance. *IRE Journals*, 2(9), pp.468-482.
- [104] Okonkwo, C.S., Agbabiaka, J., Ogunwole, O., Mayo, W. & Okeke, O.T., 2020. Model for Demurrage Elimination and Port Logistics Efficiency in Emerging Economies. *International Journal of Multidisciplinary Research and Growth Evaluation*, 1(5), pp.552-562.
- [105] Marmot, M., 2010. *Fair Society, Healthy Lives: The Marmot Review*. UCL, London.
- [106] Marmot, M. et al., 2020. *Health Equity in England: The Marmot Review 10 Years On*. Institute of Health Equity, London.
- [107] Commission on Social Determinants of Health, 2008. *Closing the Gap in a Generation*. WHO, Geneva.
- [108] World Health Organization, 2010. *A Conceptual Framework for Action on the Social Determinants of Health*. WHO, Geneva.
- [109] Victora, C.G. et al., 2003. Applying an equity lens to child health and mortality. *Lancet*, 362(9379), pp.233-241.
- [110] National Population Commission (NPC) Nigeria and ICF International, 2018. *Nigeria Demographic and Health Survey 2018*. NPC and ICF International, Abuja.
- [111] Lumley, T., 2010. *Complex Surveys: A Guide to Analysis Using R*. Wiley, Hoboken.
- [112] Ginsberg, J. et al., 2009. Detecting influenza epidemics using search engine query data. *Nature*, 457(7232), pp.1012-1014.
- [113] Signorini, A., Segre, A.M. & Polgreen, P.M., 2011. The use of Twitter to track levels of disease activity and public concern. *PLOS ONE*, 6(5), e19467.
- [114] De Choudhury, M. et al., 2013. Predicting depression via social media. *Proceedings ICWSM 2013*, pp.128-137.
- [115] Dunn, A.G. et al., 2015. Associations between negative opinions about HPV vaccines on social media and vaccination rates. *Journal of Medical Internet Research*, 17(6), e144.
- [116] Mikolov, T. et al., 2013. Distributed representations of words and phrases and their compositionality. *Advances in NIPS*, 26, pp.3111-3119.
- [117] Socher, R. et al., 2013. Recursive deep models for semantic compositionality over a sentiment treebank. *Proceedings EMNLP 2013*, pp.1631-1642.
- [118] Nguyen, D.Q. et al., 2020. BERTweet: A pre-trained language model for English tweets. *Proceedings EMNLP 2020*, pp.9-14.
- [119] Hyndman, R.J. & Khandakar, Y., 2008. Automatic time series forecasting: The forecast package for R. *Journal of Statistical Software*, 27(3), pp.1-22.

- [120] Box, G.E.P. & Jenkins, G.M., 1970. *Time Series Analysis: Forecasting and Control*. Holden-Day, San Francisco.
- [121] Scott, S.L. & Varian, H.R., 2014. Predicting the present with Bayesian structural time series. *International Journal of Mathematical Modelling and Numerical Optimisation*, 5(1-2), pp.4-23.
- [122] Zhang, G.P., 2003. Time series forecasting using a hybrid ARIMA and neural network model. *Neurocomputing*, 50, pp.159-175.
- [123] Nsoesie, E.O. et al., 2014. A systematic review of studies on forecasting the dynamics of influenza outbreaks. *Influenza and Other Respiratory Viruses*, 8(3), pp.309-316.
- [124] Aboagye-Sarfo, P. et al., 2015. Modelling and forecasting emergency department demand in Western Australia. *Journal of Biomedical Informatics*, 57, pp.62-73.
- [125] Wirtz, V.J. et al., 2017. Essential medicines for universal health coverage. *Lancet*, 389(10067), pp.403-476.
- [126] Michael, O.N. & Ogunsola, O.E., 2019. Determinants of Access to Agribusiness Finance and Their Influence on Enterprise Growth in Rural Communities. *Iconic Research and Engineering Journals*, 2(12), pp.533-548.
- [127] Michael, O.N. & Ogunsola, O.E., 2019. Strengthening Agribusiness Education and Entrepreneurial Competencies for Sustainable Youth Employment in Sub-Saharan Africa. *Iconic Research and Engineering Journals*, 2(9), pp.416-431.
- [128] Lilian, I.N., Liadi, K.O., Yeboah, T.J. & Apelehin, A.A., 2020. Understanding Cross-Cultural Communication: Identity, Diversity, and Global Interaction. *Gyanshauryam, International Scientific Refereed Research Journal*, 3(4). <https://doi.org/10.32628/GISRRJ21352>
- [129] Dagodzo, D., 2018. A Conceptual Framework for UAV Integration into National Power Grid Inspection Programs. *IRE Journals*, 2(5), pp.391-412.
- [130] Dagodzo, D., 2018. A Review of UAV Applications in Electrical Transmission Line Inspection: Methods, Technologies, and Challenges. *IRE Journals*, 2(6), pp.234-254.
- [131] Dagodzo, D. & Ahiaeke Patrick, M.C., 2020. UAV-Based Pipeline and Corridor Monitoring: A Review of Current Practices and Emerging Technologies. *IRE Journals*, 3(10), pp.574-597.
- [132] Eyetsemitan, R.A., Ambali, K.B., Oyeleye, A.O. & Fadayomi, O., 2020. Multi-Stakeholder Governance Alignment in Joint Venture Operations: A Conceptual Framework for Coordinating Business Processes in Highly Regulated Environments. *IRE Journals*, 4(4), pp.418-441.
- [133] Aye, P.A. & Tawose, O.M., 2016. Physiological responses of West African dwarf sheep fed graded levels of *Gmelina arborea* leaf and cassava peel concentrates under different management systems. *Agriculture and Biology Journal of North America*, 7(4), pp.185-195.
- [134] Aye, P.A. & Tawose, O.M., 2015. Acceptability and utilization of graded levels of *Gmelina arborea* leaves and cassava peels concentrate by West African dwarf sheep. *International Journal of Advances in Agriculture*, 4(2), pp.415-422.
- [135] Chawla, N.V. et al., 2002. SMOTE: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16, pp.321-357.
- [136] He, H. & Garcia, E.A., 2009. Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering*, 21(9), pp.1263-1284.
- [137] Fernandez, A. et al., 2018. *Learning from Imbalanced Data Sets*. Springer, Berlin.
- [138] Chicco, D. & Jurman, G., 2020. The advantages of the Matthews correlation coefficient over F1 score and accuracy. *BMC Genomics*, 21(1), pp.1-13.
- [139] Fawcett, T., 2006. An introduction to ROC analysis. *Pattern Recognition Letters*, 27(8), pp.861-874.

- [140] Bradley, A.P., 1997. The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Pattern Recognition*, 30(7), pp.1145-1159.
- [141] Mbonu, I.S., Aliliele, C., Iwuanyanwu, U. & Uzoka, E., 2020. A Review of Identity and Access Management Integration Strategies in Hybrid and Multi Cloud Environments. *International Journal of Multidisciplinary Research and Growth Evaluation*, 1(5), pp.795-810.
- [142] Mbonu, I.S., Aliliele, C., Uzoka, E. & Oluoha, O.M., 2019. A Review of Comparative Data Protection Regulations and Secure Cloud Implementation Strategies Across Jurisdictions. *Iconic Research and Engineering Journals*, 2(9), pp.482-501.
- [143] Mbonu, I.S., Aliliele, C., Iwuanyanwu, U. & Oluoha, O.M., 2018. A Conceptual Framework for Legal and Ethical Risk Modeling in Enterprise Data Protection Governance Systems. *Iconic Research and Engineering Journals*, 2(2), pp.207-226.
- [144] Mbonu, I.S., Iwuanyanwu, U., Aliliele, C. & Uzoka, E., 2020. Advances in Infrastructure as Code Governance for Secure Terraform Based Enterprise Cloud Deployments. *International Journal of Multidisciplinary Research and Growth Evaluation*, 1(5), pp.811-828.
- [145] Mullen, C., 2020. *Snowflake for Dummies*. Wiley, Hoboken.
- [146] Sanni, J.O., Ajiga, D. & Atima, M.E., 2020. Analytical Models Addressing Measurement Challenges of Marketing Return on Investment in Regulated Services. *International Journal of Multidisciplinary Research and Growth Evaluation*, 1(5), pp.636-648.
- [147] Sanni, J.O., Ajiga, D. & Atima, M.E., 2020. Systematic Review of Product Management Strategies in Mobile Network Rollouts Across Emerging Markets. *International Journal of Multidisciplinary Research and Growth Evaluation*, 1(5), pp.661-673.
- [148] Arumosoye, O.M. & Obriki, O.D., 2020. A Governance-Oriented Conceptual Model for Contractor Safety Performance in Multi-Contract Industrial Projects. *International Journal of Multidisciplinary Research and Growth Evaluation*, 1(5), pp.728-740.
- [149] Arumosoye, O.M. & Obriki, O.D., 2018. Development of an Integrated Heat Stress Risk Conceptual Model for Industrial Operations in Extreme Environments. *IRE Journals*, 1(12), pp.141-160.
- [150] Arumosoye, O.M. & Obriki, O.D., 2019. Systematic Review of Near-Miss and Hazard Observation Data Utilization in Industrial Safety Management. *IRE Journals*, 3(2), pp.981-999.
- [151] Obriki, O.D. & Arumosoye, O.M., 2019. A Conceptual Framework Linking Management Safety Walkthrough Frequency and Coverage to Safety Culture Outcomes in Mega Projects. *IRE Journals*, 2(8), pp.355-374.
- [152] Obriki, O.D. & Arumosoye, O.M., 2018. Conceptual Modeling of Data-Driven Occupational Safety Risk Control in Large-Scale Energy Infrastructure Projects. *IRE Journals*, 1(7), pp.169-189.
- [153] Obogo, S.F., Arumosoye, O.M. & Obriki, O.D., 2020. Advances in Internal QHSE Audit Systems for Industrial Engineering Operations. *Iconic Research and Engineering Journals*, 4(4), pp.399-417.
- [154] Obogo, S.F., Arumosoye, O.M. & Obriki, O.D., 2020. Conceptual Risk Management Model for Heavy Lifting and Crane Installation Engineering Operations. *Shodhshauryam, International Scientific Refereed Research Journal*, 3(4), pp.122-144.
- [155] Obogo, S.F., Arumosoye, O.M. & Obriki, O.D., 2020. Critical Review of Occupational Safety Management Systems in Oil and Gas Maintenance Projects. *Shodhshauryam, International Scientific Refereed Research Journal*, 3(4), pp.145-166.
- [156] Mbonu, I.S., Iwuanyanwu, U., Uzoka, E. & Oluoha, O.M., 2019. Advances in Enterprise Log Analytics and Automated Incident Response Architectures Using Python and SIEM Platforms. *Iconic Research and Engineering Journals*, 3(2), pp.1000-1019.

- [157] Ioffe, S. & Szegedy, C., 2015. Batch normalization: Accelerating deep network training. Proceedings ICML 2015, pp.448-456.
- [158] Mitchell, T.M., 1997. Machine Learning. McGraw-Hill, New York.