# Artificial Intelligence in Automated Tax Auditing and Risk Scoring

## Madhu Sathiri

Independent Researcher, India

**Abstract:** Tax compliance constitutes a substantial challenge for national revenues and public services worldwide, particularly in a digital economy that enables rapid international transactions. Artificial intelligence (AI) can enhance automated risk scoring and tax auditing capabilities by bridging the gap between the rapid development of machine-learning methods and the pressing operational needs of tax administrations. The applicability of AI-based risk scoring and auditing methods in the tax domain has thus far remained largely unexplored in the literature, as has the evaluation and validation of the resulting systems. Motivation, design, methods, and specific foundations (data-driven evidence, risk-scoring models, and automated auditing techniques) are presented in these sections, along with considerations of data governance, privacy, and ethics.

Evidence drawn from knowledge engineering and computational taxonomy outlines the data requirements, provenance, and quality for reliable AI applications for tax compliance, providing a foundation for subsequent sections on risk-scoring models, data-driven evidence, and automated tax auditing. Risk-scoring models identify the relevance of explainability, novelty detection, and machine-generated human-readable components, supported by privacy-preserving techniques and algorithmic transparency. Two key approaches are identified: supervised learning generates predictions for tax-relevant domains, whereas unsupervised and semi-supervised methods support hierarchical anomaly detection. These directions together address the completeness of AI auditing systems, complementing research on planning, knowledge representation, and evaluation of audit systems.

**Keywords:** Automated Tax Auditing. Artificial Intelligence; Classification and Regression; Data-Driven Audit Planning; Data Mining Technologies; Document Analysis; Natural Language Processing; Risk Scoring Models. Auditing Apparatus. Governance Framework.

## 1. INTRODUCTION

Governance of the vast digital economy relies on tax compliance, and fraud and evasion cost governments hundreds of billions of dollars annually. Cutting-edge technologies can help tax authorities maintain confidence in the integrity of tax systems. Applying artificial intelligence as an umbrella term for powerful data-driven, statistical, computational, and predictive technologies for dealing with real-world and often real-time problems, systems and processes can be developed to automate risk scoring for the likelihood of a tax audit and automated systems to help digitise tax audit processes.

Tax authorities collect vast amounts of data from multiple sources, with substantial investments in data science teams, yet only a small proportion of taxpayers is audited each year. The overall cost-benefit ratio of auditing potentially fraudulent tax returns is low, as most are likely to be compliant. Accurate risk classification of taxpayers and prediction of audit recommendations, including if a business is required for audit and at what state of the audit cycle, can improve effectiveness and efficiency even with no increase in the number of auditors. By concentrating scarce resources on high-risk tax returns, these systems can improve effectiveness and efficiency even with no increase in the number of auditors. These innovations will provide future research avenues and the foundations for the next generation of automated decision systems for dealing with real-world problems.



Fig 1: AI in Tax

## 1.1. Background and Significance

Tax compliance remains a persistent challenge in countries operating digital economies. Despite improvements in technology and rising user acceptance of e-government services, taxpayers do not complete lodging their returns correctly or, in some cases, at all. To mitigate this, tax authorities deploy tax risk-scoring and audit-planning models that predict where compliance resources will be best allocated. A growing number of tax authorities are pursuing AI solutions to support these models. A collection of theoretical building blocks is presented to support this emerging direction of research.

Compliance phenomena must be operationalized for AI to be of relevance. A complete inventory of compliance-related phenomena must therefore be built up, replete with all the requisite details, and a taxonomy structured around the nature of the phenomenon would permit cyberspace data to be allocated to the appropriate areas. The extent to which AI is able to assist with the discovery and design of risk-scoring models, the computation of risk scores, and the plan for risk-based tax audits can then be established. The potential of AI is that these activities can be automated with minimal human intervention.

**Equation 1: Precision, Recall, TPR, FPR (metrics derived from the counts)**
From the confusion matrix:
**Precision**

$$\text{Precision}(\tau) = \frac{TP}{TP + FP}$$

**Recall / True Positive Rate (TPR)**

$$\text{Recall}(\tau) = \frac{TP}{TP + FN} = \text{TPR}(\tau)$$

**False Positive Rate (FPR)**

$$\text{FPR}(\tau) = \frac{FP}{FP + TN}$$

**F1 score**

$$F1(\tau) = \frac{2 \cdot \text{Precision}(\tau) \cdot \text{Recall}(\tau)}{\text{Precision}(\tau) + \text{Recall}(\tau)}$$

## 1.2. Research design

Compliance with tax obligations is essential for the sustainability of welfare states. Pressure is therefore mounting on public authorities to step up audits of individuals' and companies' tax returns and payments. Digital economies pose unprecedented challenges to controlling tax compliance, but also offer access to vast amounts of data. AI techniques help squeeze sense from these masses of data and detect compliance breaches. Four tax-auditing areas particularly benefit from AI: computing risk scores, automated audit planning, detecting anomalies, and examining documents. In all these tasks, performance and quality assurance must consider failings, hazards, and limitations peculiar to each AI-type technique. Data management safeguards need to ensure that the data used to both train and operate AI systems meet stringent quality criteria. Built-in privacy precautions are paramount, and adherence to the concept of privacy-by-design is a key prerequisite. AI systems demand the highest standard of transparency and accountability.

Research was based on current literature and expert interviews. The aim was to delineate the potential of AI for automating tax auditing and risk scoring as well as the related challenges and dangers. Insights garnered were directed at guiding future methodological development. The resulting high-level categorization of AI opportunities provides a useful roadmap while simultaneously spotlighting the main pitfalls, and associated preventative or mitigating measures. The AI opportunities emerge from understanding the inadequacies of human auditors and the potential of techniques based on supervised machine learning, unsupervised or semi-supervised learning, and natural language processing.

## 2. THEORETICAL FOUNDATIONS OF AI IN TAX COMPLIANCE

The paper identifies and defines core compliance phenomena using data-driven evidence and a computational taxonomy. Such a foundation supports the development, evaluation, validation, and adoption of risk-scoring models employed in automated tax audit planning. Risk scores indicate the compliance of entities or their behaviours. Calibration measures the correspondence between scores and levels of strength or weakness. Discrimination quantifies the ability of a model, or model threshold, to distinguish among levels. For models designed to detect noncompliance, thresholds control the rate of false negatives. Scoring frameworks, calibration, discrimination, thresholding, and evaluation criteria such as area under the receiver operating characteristic curve and precision–recall curves are steps in risk-scoring processes.

Tax compliance comprises a variety of phenomena, and the data underlying compliance-related activities differ accordingly. Detecting, evaluating, modelling, and addressing these behaviours rely on suitable data-science methods

tailored to their respective data types and risk-scoring requirements. A computational taxonomy of AI in tax compliance maps compliance-related activities to data types, sources, and structures. Supervised learning techniques are evident in automated tax auditing, particularly in audit-planning systems that classify, score, or rank entities according to their risk of noncompliance. Other methods, such as risk-scoring calibration, unsupervised anomaly detection, and semi-supervised handwriting recognition, are instrumental for AI-enabled tax compliance activities.
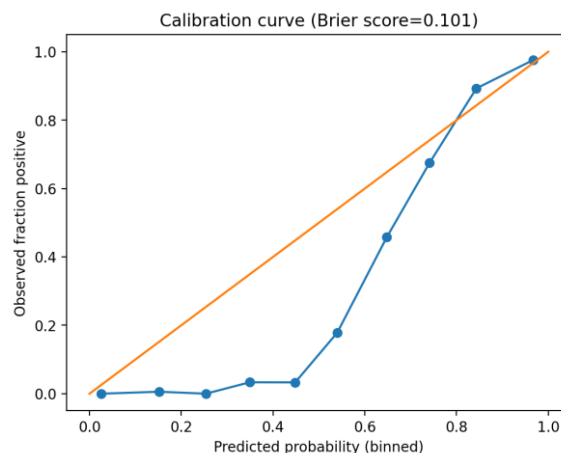
**2.1. Data-Driven Evidence and Computational Taxonomy** Compliance with tax law can be seen as a multivariate phenomenon in a rich data ecosystem characterized by highly diverse sources, types, and nature of data, and it can benefit from principled computation. Systematic, theory-guided knowledge synthesis on the comprehensive data landscape therein is therefore essential, as is the careful and computationally sound exploitation of such data landscape. By following a data-driven approach based on the descriptions provided by agency staff, the data landscape with respect to tax compliance is mapped and a taxonomic classification of compliance-related phenomena suitable for risk-scoring purposes is developed. A rich data and knowledge ecosystem provides opportunities for multivariate exploration and offers the groundwork for data-driven feature engineering.

Multivariate statistical analyses allow for the overall data landscape in a tax authority to be utilized to generate evidence for risk scorings. The results from these statistical analyses at the data-variables level then provide critical inputs and guidance for the computation of risk scores. The data-driven analysis distinguishes different families/taxonomies of compliance phenomena, namely those being commonly adopted by tax authorities; those observable from claim data; those observable from other taxpayers' behaviour; those pertaining to compliance documentation; those linked to specific law provisions; and those observable only through tax audits.

**2.2. Risk Scoring Models: Principles and Metrics**

Numerous methods, approaches, and frameworks are found in the literature for developing risk-scoring models and systems; however, two aspects are paramount: risk scoring measures (considering a generic risk-scoring study) and the calibration strategy. The basic level of risk scores is a single value for each entity—a taxpayer or a group of taxpayers—that represents the total risk or a material facet of it. Risk-scoring models differ according to these applications; that is, how the generated scores are then used. Risk-scoring models operate on the basis of the risk score lies either above or below a clearly defined predefined threshold. Explicit criteria specify what makes a risk score above average or below acceptable values.

Calibration accuracy is essential in risk-scoring studies: the percentage of "true positives" and "false positives" (or other appropriate combinations) directs the design of the models and their real-life application. Calibration testing employs a separate dataset—it must never be the validation or test dataset for establishing the "quality" of the risk-scoring model. Scoring models are developed within testing, the target data used must also have the proper level of predictions for practical application. The taxonomy uses common risk attributes across countries. Normalisation of these different sources allows the establishment of a risk-scoring framework from an extensive scale, unbiased dataset using commonly accepted methods for discrimination and calibration evaluation metrics.



Calibration curve (Brier score=0.101)

# 3. AUTOMATED TAX AUDITING: METHODS AND ARCHITECTURES

Tax-related matters and substantial areas of tax auditing can be conducted through the systematic application of AI techniques such as supervised, unsupervised, and semi-supervised learning; natural language processing (NLP); and text generation. More specifically, these implementations can be categorised as supervised learning in audit planning,

unsupervised and semi-supervised techniques for anomaly detection, and NLP for document analysis. Each category requires carefully crafted architectural considerations, and the effectiveness of these methods depends heavily on a range of practical aspects such as conceptual model design, integration into operational settings, and the availability of adequate data in useful form.

Supervised learning is commonly used for risk scoring, benchmarking, and resource allocation in tax audits; supervised anomaly detection; labelling data for unsupervised algorithms; and automating textual data, including audit report generation. In relation to audit planning, either classification or regression models can be applied to these tasks, depending on the nature of the label. For classification, each taxpayer group can be assigned a score, with low-scoring classes receiving the least attention. Regression predicts scoring as a budgetary control mechanism. Deployment also needs careful planning, with supervision essential for results that may support legal penalties or defence against appeals.
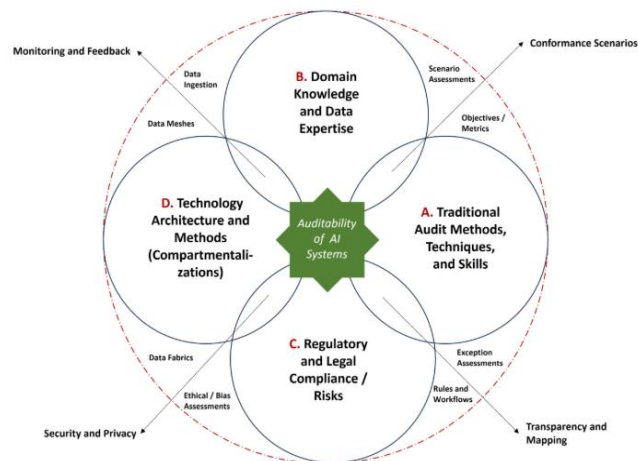


Fig 2: Automated Tax Auditing

### 3.1. Supervised Learning in Audit Planning

Tax administrations require accurate risk classification of filing cohorts, risk-factored scoring of individual filers, and detection of reporting nonconformities. Filers are groups of return cases sharing one or more operational dimensions with distinct risk characteristics. By comparing these operational dimensions—size of assets declared, type of activity conducted, and net income—across tax years, cohorts likely to deviate significantly in their reporting between years can be identified and monitored. Classification models produce a risk score over one or several phenomena. These indicate audit type (yes/no), return verification system (full selective/verify at least 50% of higher revenues/take into account all reported disallowable tax expenses), and specific caseathon triggers (yes/no).

Three modeling categories facilitate tax auditing planning: those predicting deviation in measurable performance from year to year, identification of fl file—classified as low-value companies with no materiality concerns—and measures predicting risk of high fiscal impact. Any supervised learning model utilizes a pipeline comprising data preparation for training, model training and hyper-parameter tuning, and model application at scale. Initial exploratory analyses reveal unreported income or unreported activities. Supervised learning is used when suitable labels exist for historical and non-historical records. Probabilities are continuous-valued predictors; model validation clarifies best threshold values. Decision trees and ensembles are natural formulations for class variables; probability values are treated like scalars in regression models.

### 3.2. Unsupervised and Semi-Supervised Techniques for Anomaly Detection

Unsupervised and semi-supervised techniques can augment tax compliance monitoring efforts by uncovering evidence of illicit behaviours that would otherwise remain hidden. Clustering algorithms enable exploratory data analysis to identify non-tax-compliant behaviours that deviate from the norm, while anomaly detection methods automatically flag unusual transactions for further investigation. Semi-supervised learning that combines normal and abnormal data provides a robust foundation in cases where labelled examples of illicit behaviours are scarce or do not exist.

Unsupervised anomaly detection determines whether an object is consistent with a given dataset containing only normal examples. With an adequate representation of the natural state, changes in that state—such as fraud attempts or serious system faults—are detected even if labelled examples of the changes do not exist. However, most unsupervised algorithms depend only on the training data and are thus vulnerable to being fooled by an overly informative model. They also require all features to be purely numerical ($\pm1$ to indicate a false element). Semi-supervised methods lift these restrictions by also making use of known abnormal samples in conjunction with the (often extensive) sets of normal samples.

**Equation 2: ROC curve and AUC (discrimination)**

As you sweep $\tau$ from 1 down to 0, you get pairs:

$$\big(\mathrm{FPR}(\tau), \mathrm{TPR}(\tau)\big)$$

Plotting these gives the **ROC curve**.

**AUC** is the area under ROC, usually computed numerically (trapezoids):

$$\mathrm{AUC} \approx \sum_k \left(\mathrm{FPR}_{k+1} - \mathrm{FPR}_k\right) \cdot \frac{\mathrm{TPR}_{k+1} + \mathrm{TPR}_k}{2}$$

**3.3. Natural Language Processing for Document Analysis** Natural language processing (NLP) subsumes a variety of computational linguistics methods for text mining. In tax compliance, NLP tasks typically include information extraction from documents and multilingual information retrieval, entity recognition, and automatic summarization for knowledge distillation. NLP for information extraction and recognition relies on labelled corpora for supervised training. Yet labour-intensive annotation incurs financial costs and time lags, especially for low-resource languages such as Catalan. Hence semi-supervised and unsupervised methods may be more suitable. For automatic summarization, techniques based on extractive document summarization are often deployed owing to their robustness and data-efficient training pipelines.

The data processing involved is determined by the task at hand. For information extraction from form-like documents, template matching performed by rules or a regular expression engine constitutes a basic yet effective approach. When entity detection focuses on lists of structured information such as business registration data or VAT registration numbers, regular-expression methods also suffice. Information extraction from scanned documents requires a preceding optical-character-recognition step. For named-entity recognition, an initial step typically involves the automatic production of labelled documents by combining a list of entities and rule-based text extraction. Retrieving documents on specific topics from the Internet requires the definition of the topic and its translation into the languages of interest. For summarization, document-cleansing procedures such as sentence-splitting and tokenization precede the computation of sentence importance scores. Evaluation relies on automatic metrics such as ROUGE and BERTScore.

## 4. Data Governance, Privacy, and Ethics

The use of data, especially data about individuals, raises privacy concerns. Privacy by design means that legislation must be complemented by privacy-centric design principles such as data minimization (retaining only the attributes needed for business processes), purpose limitation (using data only for the initial purpose), and incorporating privacy safeguards throughout data life cycles. Furthermore, access to data must be restricted to required users. Data privacy must also comply with existing law and regulation. The family of General Data Protection Regulation in the European Union and equivalent legislations elsewhere provides such a framework. In addition, compliance operates at two levels: the legal obligations that naturally emanate from the legislation and the design principles that facilitate compliance in practice.

The results of machine learning and data management processes need to be understandable not just to developers, data scientists and auditors, but also to the end-users and the stakeholders concerned. Explainability and interpretability of results increase the value of tax administration activities, contributing to the objective of the administration to increase voluntary compliance. In addition, machine learning processes generate system features that aid transparency and accountability. Audit trails that record the machine learning process for model specification and operationalization facilitate result interpretation, while explanations for model results help determine actions to be taken where harm is possible. However, explanation is not the only requirement; information that audits were developed adds credibility to decisions such as granting licenses.

The deployment of artificial intelligence in tax-related processes considered by multiple governmental entities has created public concerns about possible bias in training data, adverse impact against groups protected by nondiscrimination legislation, and negative effect on special-interest groups. Addressing these concerns involves technical assessment of biases, fairness indicators, and expert stress-testing of supervised models, as well as prospective evaluation of classes at risk with appropriate mitigating strategies. Nevertheless, skillful use of the artificial intelligence-enhanced processes can still create net economic benefits for all parties including affected groups. For instance, voluntary use of advance personal-expenditure tax credits can be structured to benefit the least and middle income classes while minimizing harm to the higher income classes.

### 4.1. Data Quality and Provenance

Data processes must define how each datum is created, transformed, and originally sourced. Information lineage is essential for understanding risk and for compliance purposes. Sufficient data quality is a prerequisite for trust in methods employing such data, since errors or biases can lead to spurious conclusions. Data quality is multi-dimensional, and these
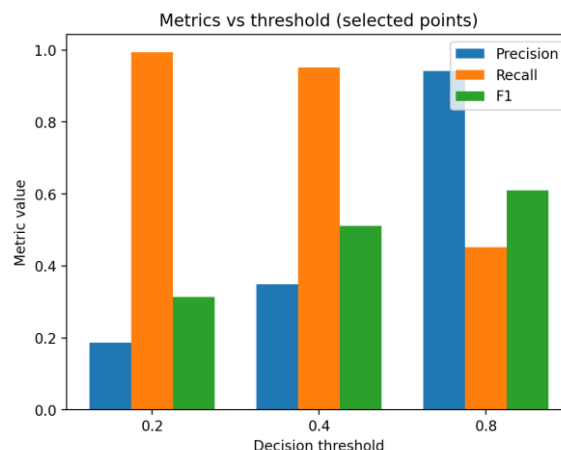
dimensions require validation checks capable of exposing flaws. Provenance metadata helps guarantee reproducibility— the ability of an independent auditor to re-create earlier analyses and confirm (or otherwise) their results.

Provenance metadata for two AI systems is summarized by the Plan Information Model in Big Data. Its data quality dimensions are confusing, since data validity, accuracy, and credibility are redundant. The conventional list comprises dimensions such as accuracy (free from error), completeness (free from missing pieces, both at the level of records and at the individual attributes), uniqueness (no duplicate records), consistency (no conflicting values), and timeliness (fit-for-purpose temporal properties). Metadata should check for these ontological and documented aspects. Validation mechanisms include spam filters, grammar-checking heuristics, and implementations of the previously mentioned techniques for anonymity detection and mitigation. A few data sources, such as the volume of data generated, already provide such checks and associated feedback when data quality fails, thereby affecting their practical utility.

### 4.2. Privacy By Design and Regulatory Compliance

Privacy safeguards must be embedded in the AI solutions and services used in tax auditing systems and processes by enforcing appropriate technical and organizational controls. Tax auditing operations should comply with applicable laws and regulations related to data protection, information security, and privacy. Such regulations usually require that consumer data is handled in a secure manner and shared with third parties only in exceptional situations.

Privacy By Design encompasses the entire lifecycle of a project—from inception to completion and beyond. AI systems should collect and retain only the information that will be used for the analysis. Data that is not being used should be periodically deleted. Access to personal data should be minimized and limited to those individuals whose job responsibilities entail such access. High-risk data such as directly identifiable data should be stored in encrypted databases. Data access controls should be in place to prevent unauthorized access to personal data. Data exchange or sharing for other purposes should be done only in circumstances mandated by the law or in emergency cases such as life-threatening situations.



Metrics vs threshold (selected points)

### 4.3. Algorithmic Transparency and Accountability
Systematic explanations of AI decision-making enhance trust and facilitate oversight. Rationale articulation in natural language may foster acceptance among agents affected by key decisions. Explanation requirements can be codified as restrictions on the solutions of supervised learning problems. Discrimination by attributes protected under human rights constitutes an ethical bias that can be filtered using fairness constraints.

No matter how complex a predictive model, the ability to reconstruct past decisions allows affected parties to question their appropriateness and their underlying reasoning. The audit trail following a negative Wood Score can include a suitable summary of the detection data and a justification of the model output. Appropriately trained InfoDumpers can provide an overview of the detection and risk profiling. They can also generate explanations for the prediction of the Scammer Model and Rumor Model and summarize the logic behind false negatives in the Detection Model. Such summaries may even appear as articles in a news feed.

Audit trails increase actionability and accountability. "Algorithmic accountability" implies supervision by specialized bodies capable of understanding the algorithm and its rationales. The methodology proposed in provides the foundation for such oversight. Developing transparent and accountable AI tools for tax compliance and audit, and enabling access to the common logic present in all functions underpinning the Monitoring Model, would create a new paradigm for human–AI collaboration: a human-in-the-loop framework in which AI acts as an enhancer, augmenting rather than replacing professionals. The contribution and rationale behind critical decisions can be retrieved by auditing functions. Emphasizing the AI-enhancer role also helps address privacy concerns and creates a trustworthy ecosystem for dialogue with the community of agents directly affected by the decisions of the tax administration.

## 5. EVALUATION AND VALIDATION OF AI SYSTEMS IN TAX AUDITING

Determining whether a solution works as intended and delivers tangible benefits can be difficult. Consequently, AIs designed for use in tax administration should incorporate thoroughly tested methods that demonstrate their effectiveness ex ante and ex post.

Performance Metrics and Benchmarking

Performance indicators are crucial for evaluating success and determining whether a deployment is desirable. Baselines for comparison should be established early in development; these may come from exploratory data analysis, earlier research, domain knowledge, or other sources. Quantitative metrics are essential for supervised learning systems, and while absolute values may suffice for choice problems, they should always be interpreted relative to baseline alternatives for ranking selections. Augmenting samples with synthetic data allows additional assessments, such as reinforcement learning evaluations when simulating an agent operating within a custom environment. Cross-validation techniques are needed in resource-constrained applications, while third-party datasets enable independent benchmarking of supervised methods.

Bias, Fairness, and Robustness Assessment

Where risk scoring systems produce sensitive attributes, bias-detecting statistics can quantify potential unfairness. Resilience to boundary-case data, domain shifts, adversarial perturbations, and other stress tests reveals susceptibility to predictable errors. Mitigation strategies include design principles. Ensemble methods incorporate heterogenous depth, width, or modelling techniques to provide variance-robustness; adversarial-dropout-like approaches curate differentially-stable groups. Robustness-promoting labels identify safe regions and abundant climates. Pre-processing techniques, such as re-weighting, counterfactual-data generation, and adversarial practices, reduce bias. After-care methods, including calibration and re-weighted fairness optimizations, rectify biases post hoc.

Audit Effectiveness and Return on Investment

Effectiveness metrics assess how well predictions align with actual outcomes for autonomous agents. Cost-benefit analyses indicate economic viability during deployment, disclosing costs of additional actions. Studies into audit approaches and systems that incorporate worldwide data divulge their expected value. In nature, public goods remain stable. Components that incur excessive expenses require extensive training time; hence, latent demand—anticipated lowered tax collection via system prediction—indicates appetites. Exploratory data analysis determines the effectiveness of prediction-type methods. Strategic value, driven by skilled employee shortage and recruitment issues, endorses training and reallocation beyond tax behaviour discovery.
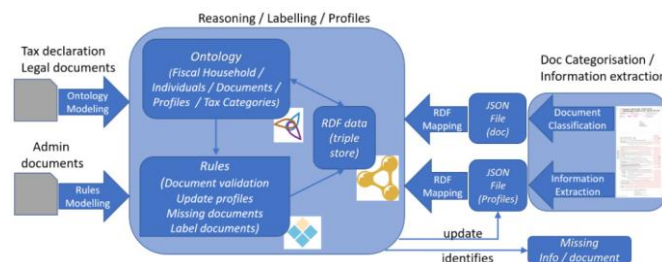


Fig 3: Evaluation and Validation of AI Systems in Tax Auditing

### 5.1. Performance Metrics and Benchmarking

Success criteria depend on the AI type and system purpose. For supervised learning in audit planning, success signifies predictive accuracy for strategies. For risk scoring models, a well-calibrated, discriminative score–threshold combination enhances compliance management and tax authority reputation. For unsupervised anomaly detection, successful systems expose clustered anomalies. The entire auditing domain can be supported by labelled/natural language processing (NLP) tasks. Cross-validation schemes for supervised methods ensure performance generalizability. Inherent biases and fairness concerns require careful handling. Biased decisions harm taxpayers and undermine tax authorities. Conformity with legal and institutional standards, stakeholder disadvantage mitigation, and service fairness are vital. Stress testing assesses sensitivity to various factors, including data shifts and training scope. Cost-benefit analysis quantifies monetary gains relative to effort and investment. Strategic business value describes qualitative effects on authority reputation and ecosystem integrity.

Task success and benchmarking metrics enable objective monitoring. Classification/regression tasks employ traditional success measures and proxy scores. Ranking-based model performance uses the area under the receiver operating characteristic curve (AUC) and precision–recall measures. Novelty detection tests distinguish button-type from multi-alyze multi-class divergence discrimination. Relational similarity implements standard and cross-domain similarity for taxonomy-driven evaluation. NLP task validation employs conventional metrics and human assessment.

## 5.2. Bias, Fairness, and Robustness Assessment

When taxpayer compliance scores are based on personal data, algorithmic bias and unfairness can lead to disproportionate surveillance and sanctions for highly scrutinized groups. Bias stems from biased training data and algorithms that amplify correlations—notably, a tendency of male taxpayers to be less compliant, at least in some jurisdictions. Thus, for personal data-driven estimates of compliance and related measures of fairness, bias in data, algorithm, and outcome should be assessed with fairness metrics and stress tests. For example, the bias mitigation literature in computing provides a wide variety of strategies, including preprocessing of training data (e.g., oversampling, exclusion of features correlated with bias), in-processing bias correction (e.g., adversarial training), and post-deployment adjustment.

Robustness refers to the ability of the AI-powered applications to maintain their performance when facing data shifts. Addressing robustness usually entails stress testing applications against foreseeable shifts in input distribution or model parameters (e.g., sub-optimal hyperparameters in the supervised learning pipelines). Beyond the recommend robustness checks against data shift, parameter uncertainty, and adversarial perturbations, the practical deployment architecture should also insulate applications against shifts in data or label distribution, which may not be apparent at the development stage. In particular, a decline in return on investment could instigate the introduction of a novel scheme, such as tax break for investing heavily in countries with a declining score.

## 5.3. Audit Effectiveness and Return on Investment

Measuring an automated tax auditing system's effectiveness requires comparing expected and actual tax revenue. Further, return on investment is evaluated by comparing the costs of training and operating the system to the benefits it provides. Through a combination of interviews, case studies, and user surveys, external reviewers assess whether the proposed models would improve efficiency and effectiveness if implemented within a tax agency. These factors are then considered alongside the five major criteria for developing successful AI systems. Specifically, governance and policy frameworks influence audit effectiveness, largely by establishing whether there are proper checks on the use of automated decision-making systems (such as the right of appeal).

The effectiveness of an automated tax auditing system depends to a significant extent on the context in which it operates, making return on investment difficult to calculate ex ante. Nevertheless, there are several cost components that should be considered. In particular, whether the construction of the system requires a large private-sector consultancy, and whether the privatised operation of some tax functions leads to the collection of less tax than would be the case should agencies carry out these tasks.

## 6. IMPLEMENTATION CHALLENGES AND OPERATIONAL CONSIDERATIONS

Many sources warn about the potential loss of human judgment and court's discretion in the use of AI for tax auditing, even labeling the risk as harmful. However, in standard environmental scanning, humans use judgment to avoid potential erroneous temptations of AI. The design process of AI systems focused on auditing must apply the human-in-the-loop design philosophy to minimize the potential overreliance problem.

Integration with existing or legacy systems represents another serious challenge related to the plumbing aspect of data-centric AI. Data integration enables the automation of any planned auditing process protecting the preparation, from data provenance and preparing all data registers for the AI system to integrate any external or third-party validation sources with automated pipelines.

Deployment in Tax Authorities and suitable Change Management is important to ensure user acceptance and full commitment. People normally fear change and resistance usually comes from a feeling of lack of autonomy or decreased trust in judgement by their hierarchical levels. New Ai systems should be extensively tested not only institutionally by IT professionals with knowledge in the areas related to AI technology, namely Data Engineering, Software Engineering, Islamics and Talent, but especially in the end user area. Tax Auditors' Suggestion About Objects, Properties or Algorithms Are Especially Important to Identify Desired Features and Satisfy Auditor´s Acceptability.

**Equation 3: Threshold selection tied to audit cost/benefit (simple ROI rule)**

A minimal expected-value decision rule:

Let:

- $C_{aud}$ = cost per audit

- $R$ = recoverable revenue if noncompliance is found

- $p = s$ = predicted probability of noncompliance

Audit when:

$$p \cdot \mathbb{E}[R \mid y = 1] - C_{aud} \geq 0 \quad \Rightarrow \quad p \geq \frac{C_{aud}}{\mathbb{E}[R \mid y = 1]}$$

### 6.1. System Integration and Interoperability

Effective auditing relies on high-quality and well-structured data that enables consistency and comparability of results. System integration requires a coherent architecture for both automated and human-assisted processes. An enterprise-level integration scheme can facilitate the required collaboration by coordinating data access and processing while ensuring that the proposed AI functionality is easily available to all relevant users.

The software engineering design principles of the system integration architecture build on wider enterprise architecture considerations within the audited organization. Integration at the technical level is structured around the integration of processes, systems, and data. The integration of processes groups the functions of conducting audits and risk assessment in such a manner that, notwithstanding the enhancements of automation, human supervision is justified at selected points in the process for the assurance of both effectiveness and explainability. The integration of systems entails the provision of a common data dictionary and the orchestration of different applications to guarantee that only one application at a time is executing the same function.

### 6.2. Change Management and Skill Requirements

Successful adoption of AI tools for tax compliance auditing depends not only on technical capabilities but also on organizational readiness for such transformations. In preparation for the rollout of newly developed capabilities, personnel involved in the audit work must acquire or further hone essential AI-related skills and knowledge, creating an environment that fosters acceptance and encourages appropriate use. Such preparation may involve role-specific training as well as more general sensitization workshops on privacy safeguards, security policies, and risk perception. Particular attention should also be given the needs of staff responsible for system administration, model supervision, and training custom modules.

Sensitivity to change extends beyond staff training. In completing the Data Governance, Privacy, and Ethics dimensions, direct stakeholders sought to identify human production information deficiencies and implement a practical solution to address them. Such concerns exist not only because the project is government-driven, but also because the adoption of supervised machine-learning models requires sufficient labeled training data from previous audits. With limited historical labeling undertaken by inspectors, the availability of such datasets may restrict immediate deployment or require alternative unsupervised detection models within some modules. Ensuring appropriate engagement of different tax-internal stakeholder groups involved in the detection process is also key.

### 6.3. Risk of Overreliance and Human–AI Collaboration

Cautioning against excessive reliance on AI while acknowledging its increasingly autonomous use in areas such as tax auditing underscores the importance of human oversight. Despite compelling evidence that AI typically improves detection rates in audit planning, AI outputs should be considered suggestions rather than directives to auditors, particularly given the serious implications of missed or incorrect audit calls. The lack of formal human–AI collaboration frameworks increases the likelihood of inferior outcomes through reduced human involvement, insufficient scrutiny of AI decisions, or misplaced trust in the AI model. Directives within AI models should continue to be supported by sufficient rationale. An early governance proposal for the Tax Administration of Enrichment in Santiago, Chile, endorses testing these safeguards within operational environments prior to broader deployment and enhancement.

By designing a human-in-the-loop architecture for audit planning—empowered by risk analysis, enabling software decision support, and functioning as an assistive robotics system—overreliance may be avoided. Creative and complex patterns of illicit behaviour are likely to persist, so skilled human resourcefulness will remain essential for their detection. Risks of overreliance will also be mitigated by monitoring AI performance and constructing dashboards to highlight areas where continued scrutiny is warranted.

## 7. LEGAL AND POLICY IMPLICATIONS

7.1. Compliance with Tax Law and Regulatory Standards

Each AI application for automated handling of tax compliance should be compliant with the relevant tax laws and regulations. Candidates in these significant or sensitive areas must be subject to additional compliance checks to ensure that the application of AI does not inappropriately change the requirements, scope, or conditions of the relevant regulation. Tax agencies retain full responsibility for the application of the AI-based systems, and the governing laws and regulations must clearly stipulate that these systems must not result in tax decisions with legal force without further human validation.

In addition to compliance with data protection and privacy provisions, the deployment of AI in Tax Administrations should also comply with other applicable legal and regulatory standards, such as those concerning the Law on State Statistics or public procurement law, where relevant. In the event that the AI model that is being developed gains decision-making powers, its operation must also comply with relevant requirements of the Administrative Procedure Code governing the authority and conduct of the tax administration. While it would enhance the efficiency of the tax administration, the legal provision for a decision-making model would require an additional sensitive policy discussion as it not only has a direct impact on public trust but may also hinder legal recourse.

**7.1. Compliance with Tax Law and Regulatory Standards** Tax compliance depends on the vision of the tax authorities, on how often and why they decide to audit. It also depends on the ability of the authority to define the tax rules and define information requests without requiring an authorization request. Furthermore, AI assistance should be validated to avoid future challenges or be supported by sources of evidence that go beyond proofbased law. The risk-scoring models should be audited periodically because the use of non-automated validation increases the risk of possible tax law infringements. Audit planning is the process that assesses whether an audit is justified, and if so, establishes the frequency, selection, and focus of audits. Tax compliance is the active behaviour of tax fulfillment models. AI-based risk assessment can help to determine which tax positions are a priority for audit enforcement resources. The tax authority can switch to an on-demand audit model where AI produces alert notifications, so all alerts are investigated, and some trigger real-time evidence requests based on the risk profile of the taxpayer. For the risk model to be effective, it must not increase the cost of compliance for the parties.



Fig 4: AI in Tax Law

**7.2. Governance Frameworks for AI in Public Sector Auditing**
Governance of AI applications in automated public-sector auditing should entail responsible policy development, risk-based supervision, and adequate enforcement. Primary statutory constraints warrant attention, especially those to ensure algorithmic fairness and transparency. Provisions requiring parties liable for taxes to document transactions in national currency (and third-party document creation) reduce the risk of collusion and risk scoring models' potential harm, while a formal regulatory framework on the use of AI in tax audits—including public consultation and professional assessment of underlying data—helps ensure algorithmic correctness. The weight of evidence should also inform the planning stage to minimize the auditing burden of non-targeted parties. Regular technical audits of algorithms and sufficient qualifications for auditors of AI-enable procedures further govern risk management.
Sensitivity to fairness, transparency, and practical enforceability should be part of public-interest protections by governance models. Stakeholder committees from the relevant sector should periodically assess potential algorithmic-harm risks across the implemented audit portfolio and introduce mitigation measures such as impact assessments, restoration or compensatory mechanisms, and inquiries for affected parties. A clear delineation of accountability is essential: each detection system needs one or more responsible entities for subsequent fairness assessments addressing potential discriminatory effects, noise, and model collapse. Other committees provide horizontal, sectoral, legal- and regulation-development oversight to advocate for fairness, proportionality, and legal compliance. Finally, the effectiveness of AI-supported public-sector audits should be a governance priority.

# 8. CONCLUSIONS

Methodological developments and signalling effects on compliance behaviour stand out as significant trends. Methodological advancements include increased adoption of semi-supervised and unsupervised techniques for compliance risk scoring and anomaly detection, as well as natural language processing for extracting information from textual documents. Anomaly detection, often performed in conjunction with clustering, is gaining prominence partly because tax auditors do not have the luxury of labelled data in massive quantities. The need to consider the full OTLC – including understanding the tax authority's target population, identifying potential non-compliers or anomalies before an audit, and evaluating audit results to hone risk-scoring systems for the future – is being adequately recognised. But

machine learning is at present being used to support only segments of this OTLC. Supervised techniques continue to dominate when creating tax compliance risk scores. In domestic tax systems, supervision refers to those revenue authorities that possess both the responsibility and the legal authority to tax based on evidenced income and expenses for individuals and corporations alike. Two key aspects of risk-scoring systems merit further attention: the monitoring of administrative performance and a research agenda focused on creating more optimal thresholds for risk.

Additional key topics warranting research are the likely signalling effects of tax authority communication, especially concerning the introduction of sophisticated AI and machine-learning support for tax auditing. Given the technological imperatives of informing future AI-enhanced tax compliance, the emphasis on surveillance via risk assessment and the consequential foreshadowing of likely taxpayer targeting are urgent considerations. Although discussions of AI and tax compliance often centre on revenue authorities' need to mine large taxpayers' digital footprints to identify and address likely risk events, delicate reputational issues remain. Dialogue about balancing legitimate privacy concerns against the need to prevent tax evasion, money laundering, or other criminal activity involving financial dealing with governments through de facto use of a digital-delivery model is also critical.

## 8.1. Emerging Trends

Automated audit planning, anomaly detection, and document analysis represent three core areas of AI application in tax auditing space. Supervisory learning techniques supervise the audit planning process. Supervised and unsupervised machine learning techniques are generally used to uncover anomalous behaviours and to detect outliers in clusters of taxpayers. Clustering techniques help identify groups of taxpayers with similar profiles and abnormal behaviours — such a combination of clustering and outlier detection is more efficient and powerful than the use of a single technique. Natural language processing (NLP) is applied to documents with large quantities of unstructured data, opening up new opportunities in information extraction, entity recognition and summarization.

While the application of AI in the public sector is still in its infancy, its projected benefits of improved data-driven decision-making and increased operational efficiency are compelling. Close matching of AI capabilities and services to specific requirements and also to the data-driven nature of the decisions is critical to both the success and implementation feasibility of specific initiatives. AI-enabled technology has the potential to outbreak traditional technologies for return on investment and strategic value, but careful and appropriate integration into the functioning of agencies is essential. Implementing BI-platform-like applications based on risk-scoring models significantly enhances the effectiveness of tax audits and the returns from tax audit activity while meeting tax administration objectives.

## REFERENCES

[1]     Sateesh Kumar Rongali. (2023). Explainable Artificial Intelligence (XAI) Framework for Transparent Clinical Decision Support Systems. International Journal of Medical Toxicology and Legal Medicine, 26(3 and 4), 22–31. Retrieved from https://ijmtlm.org/index.php/journal/article/view/1427.

[2]     Chan, T., Tan, C.-E., & Tagkopoulos, I. (2022). Audit lead selection and yield prediction from historical tax data using artificial neural networks. PLOS ONE, 17(11), e0278121. https://doi.org/10.1371/journal.pone.0278121.

[3]     Guntupalli, R. (2023). AI-Driven Threat Detection and Mitigation in Cloud Infrastructure: Enhancing Security through Machine Learning and Anomaly Detection. Available at SSRN 5329158.

[4]     de Roux, D., Pérez, B., M., Moreno, A., & Villamil, M. P. (2018). Tax fraud detection for under-reporting declarations using an unsupervised machine learning approach. In Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (pp. 215–222). Association for Computing Machinery. https://doi.org/10.1145/3219819.3219878

[5]     Inala, R. Revolutionizing Customer Master Data in Insurance Technology Platforms: An AI and MDM Architecture Perspective.

[5]     González Martel, C., Martín, J. L., & Rodríguez, J. (2021). Identifying business misreporting in VAT using network analysis. (Working paper).

[6] Varri, D. B. S. (2023). Advanced Threat Intelligence Modeling for Proactive Cyber Defense Systems. Available at SSRN 5774926.

[7]     Ippolito, A., & Lozano, A. C. G. (2020). Tax crime prediction with machine learning: A case study in the municipality of São Paulo. In Proceedings of the 22nd International Conference on Enterprise Information Systems (ICEIS 2020), Volume 1 (pp. 452–459). SciTePress. https://doi.org/10.5220/0009564704520459

[8 Garapati, R. S. (2023). Optimizing Energy Consumption in Smart Build-ings Through Web-Integrated AI and Cloud-Driven Control Systems.

[9]     Kleanthous, S., & Chatzis, S. P. (2020). Semi-supervised VAT fraud detection with gated mixture variational autoencoders. Knowledge-Based Systems, 188, 105048. https://doi.org/10.1016/j.knosys.2019.105048

[10]     Nagabhyru, K. C. (2023). From Data Silos to Knowledge Graphs: Architecting CrossEnterprise AI Solutions for Scalability and Trust. Available at SSRN 5697663.

[11]     Mehdiyev, N., Houy, C., Gutermuth, O., Mayer, L., & Fettke, P. (2021). Explainable artificial intelligence (XAI) supporting public administration processes – On the potential of XAI in tax audit processes. In F. Ahlemann, R. Schütte, & S. Stieglitz (Eds.), Innovation Through Information Systems (pp. 413–428). Springer. https://doi.org/10.1007/978-3-030-86790-4_28

[12]     Aitha, A. R. (2023). CloudBased Microservices Architecture for Seamless Insurance Policy Administration. International Journal of Finance (IJFIN)-ABDC Journal Quality List, 36(6), 607-632.

[13]     Vanhoeyveld, J., Martens, D., & Peeters, B. (2020). Value-added tax fraud detection with scalable anomaly detection techniques. Applied Soft Computing, 86, 105895. https://doi.org/10.1016/j.asoc.2019.105895

[14]     Keerthi Amistapuram. (2023). Privacy-Preserving Machine Learning Models for Sensitive Customer Data in Insurance Systems. Educational Administration: Theory and Practice, 29(4), 5950–5958. https://doi.org/10.53555/kuey.v29i4.10965.

[15]     Wu, R.-S., Ou, C.-S., Lin, H.-Y., Chang, S.-I., & Yen, D. C. (2012). Using data mining technique to enhance tax evasion detection performance. Expert Systems with Applications, 39(10), 8769–8777. https://doi.org/10.1016/j.eswa.2012.01.204

[16]   Nagubandi, A. R. (2023). Advanced Multi-Agent AI Systems for Autonomous Reconciliation Across Enterprise Multi-Counterparty Derivatives, Collateral, and Accounting Platforms. International Journal of Finance (IJFIN)-ABDC Journal Quality List, 36(6), 653-674.

[17]     Alexopoulos, A., Delaporta, D., Győri, A., Kotsogiannis, C., Olsson, O., & Pávková, A. (2023). A network and machine learning approach to detect value added tax fraud. (Working paper).

[18]     Gottimukkala, V. R. R. (2023). Privacy-Preserving Machine Learning Models for Transaction Monitoring in Global Banking Networks. International Journal of Finance (IJFIN)-ABDC Journal Quality List, 36(6), 633-652.

[19]   Didimo, W., Liotta, G., Montecchiani, F., & Policastro, G. (2020). A graph analytics framework for tax evasion detection. IEEE Access, 8, 86513–86525.

[20]     Kannan, S. The Convergence of AI, Machine Learning, and Neural Networks in Precision Agriculture: Generative AI as a Catalyst for Future Food Systems.

[21]     Gupta, M., & Nagadevara, V. (2007). Audit selection strategy for improving tax compliance: Application of data mining techniques. In Foundations of Risk-Based Audits: Proceedings of the 11th International Conference on e-Governance (pp. 28–30).

[22]     Sriram, H. K., & Somu, B. (2023). Next-Gen Banking Infrastructure: Designing AI-Native IT Architectures for Financial Institutions. Available at SSRN 5273819.

[23]     Basta, A., Huber, M., & Kirchgässner, G. (2014). VAT fraud detection in practice: A comparative evaluation of risk scoring approaches. (Working paper).

[24]     Komaragiri, V. B. The Role of Generative AI in Proactive Community Engagement: Developing Scalable Models for Enhancing Social Responsibility through Technological Innovations.

[25]   Alm, J., McClelland, G. H., & Schulze, W. D. (1995). Why do people pay taxes? Journal of Public Economics, 48(1), 21–38.

[26]     Chakilam, C. (2023). Next-Generation Healthcare: Merging AI, ML, and Big Data for Accelerated Disease Diagnosis and Personalized Treatment. American Online Journal of Science and Engineering (AOJSE)(ISSN: 3067-1140), 1(1).

[27]     Niu, Y. (2011). Tax audits and voluntary compliance: Evidence from firm responses following audits. (Journal article).

[28]     Challa, K. Dynamic Neural Network Architectures for Real-Time Fraud Detection in Digital Payment Systems Using Machine Learning and Generative AI.

[29]     Lin, Y., Chen, K., & Xu, X. (2015). Neural-network approaches for tax noncompliance detection: Evidence from administrative data. (Journal article).

[30]     Sai Teja Nuka, "Cloud-Based AI Systems for Real-Time Medical Imaging Analysis and Diagnostics," International Journal of Advanced Research in Computer and Communication Engineering (IJARCCE), DOI: 10.17148/IJARCCE.2022.111252.

[31]     Vellido, A., Martín-Guerrero, J. D., & Lisboa, P. J. G. (2012). Making machine learning models interpretable in public-sector risk scoring: A review with applications. (Journal article).

[32]     Challa, S. R. Next-Generation Wealth Management: A Framework for AI-Driven Financial Services Using Cloud and Data Engineering.

[33]     Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. In Advances in Neural Information Processing Systems (pp. 4765–4774).

[34]     Kiran Reddy Burugulla, J. (2023). Transforming Payment Systems Through AI And ML: A Cloud-Native Approach. Educational Administration: Theory and Practice. https://doi.org/10.53555/kuey.v29i4.10144