

# Supervised Machine Learning for Tax Evasion Detection: A Case Study with the Brazilian Tax Administration

Cleyton Pires

Universidade Federal de Santa Catarina (UFSC)

Florianópolis, Brazil

cleyton07@gmail.com

## ABSTRACT

In this study, we present an innovative approach designed to enhance the audit case selection process within the Brazilian Tax Authority (RFB) by integrating Artificial Intelligence techniques. Our emphasis is on employing supervised learning algorithms to predict the annual income of individual taxpayers coupled with outlier detection techniques to strategically prioritize cases of heightened fiscal interest. This involves leveraging a comprehensive dataset of socioeconomic variables available to the Tax Administration. A pivotal facet of our methodology is its commitment to model explainability for ensuring fairness and compliance with both legal and ethical considerations. Preliminary findings from the case study demonstrate promising results, positioning our model as a valuable complementary framework to the existing rule-based system employed in the audit case selection process.

## KEYWORDS

Tax Evasion, Supervised Learning, Outlier Detection, Explainability

### ACM Reference Format:

Cleyton Pires. 2023. Supervised Machine Learning for Tax Evasion Detection: A Case Study with the Brazilian Tax Administration. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

## 1 INTRODUCTION

Tax evasion presents formidable challenges globally, jeopardizing not only the tax base but also essential public resources designated for public goods and compromising fiscal equity. Income tax, a key revenue generator for governments, necessitates taxpayers to file annual returns, detailing income, deductions, credits, and relevant financial information for reconciling taxes withheld or paid throughout the year with the actual liability. Audits, conducted primarily by the state through tax collection agencies, constitute a crucial mechanism for ensuring tax compliance and combating evasion. Fiscal audits systematically review and verify accounting and fiscal information to ensure adherence to tax laws.

However, these audits demand significant time, human, and financial resources for both the Tax Authority and taxpayers. Given

the constraints posed by limited available resources, there is a crucial need for judicious and efficient resource utilization. This involves concentrating audit efforts on high-risk cases with substantial suspicions of tax evasion rather than allocating valuable resources to random audits or low-risk taxpayers.

The Brazilian Tax Authority (**Receita Federal do Brasil – RFB**) employs an audit case selection process grounded in a rule-based system. This process consists of two stages: (1) data cross-referencing, which aims to identify tax inconsistencies using predefined rules and conditions set by experts, resulting in a preliminary list of taxpayers generated through technical and objective criteria; and (2) individual analysis, which seeks to validate the preliminary indicators identified during the data cross-referencing stage. In essence, it involves a manual (non-automated) analysis designed to eliminate the "false positives" from the previous phase.

Despite being widely used by tax administrations, the strategy of selecting taxpayers based on pre-established rules has limitations, the main one being its limited adaptability, as rule-based systems might struggle to adapt to new, previously unseen fraud patterns. Moreover, rules need to be continuously updated and maintained as fraud tactics evolve, which can be time-consuming.

In this paper, we propose an audit case selection method based on data mining techniques to complement the traditional ruled-based system. Our specific goal is to utilize machine learning algorithms and outlier detection techniques to uncover patterns and anomalies within large datasets, using real data available to the Brazilian tax authorities. This approach aims to enhance the taxpayer selection process, particularly in terms of scalability and adaptability as machine learning models excel when there are large volumes of data to analyze and can adapt to new fraud patterns without the need for explicit programming of rules, making them well-suited for detecting evolving fraud and tax evasion tactics.

This paper is organized into five sections. Section II provides an overview of relevant literature concerning the application of Artificial Intelligence in addressing tax evasion and tax fraud. In Section III, the proposed solution is delineated, outlining its conceptual steps. Section IV presents a case study illustrating the implementation of the proposed solution, utilizing actual data from the Brazilian Tax Authority. The outcomes of this application are detailed in Section V. Lastly, Section VI outlines the conclusions drawn from the analysis of these results.

## 2 RELATED WORK

Due to the lack of availability of tax data in the public domain, the published literature on combating tax fraud and evasion is understandably scarce. Moreover, although it is a subject of significant interest, tax administrations have numerous reservations about

Permission to make digital or hard copies of all or part of this work for personal or classroom use, provided that the copies are not made for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
Conference'17, July 2017, Washington, DC, USA  
© 2023 Association for Computing Machinery.  
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00  
<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

disclosing internal projects. Due to the classified nature of taxpayer information, which requires protection by tax officers, many refrain from sharing the details of tax compliance risk projects.

In de Roux et al [3], an approach to detect property tax evasion in Bogotá, Colombia, using unsupervised learning techniques, is presented. Zumaya et al [9] utilizes electronic invoices data from the Mexican federal government to analyze interaction patterns among taxpayers and identify tax evaders, employing monthly and yearly temporal networks.

Savić et al [6] propose the HUNOD method (Hybrid UNSupervised Outlier Detection) to identify outliers in fiscal datasets, achieving an internal outlier validation rate between 90% and 98% in experiments with tax data from the Tax Administration of Serbia. Y. Lin et al [4] introduces TaxThemis, an interactive visual analysis system that assists tax auditors in identifying suspicious groups of tax evasion, especially in transactions involving related parties. It utilizes data analysis techniques and complex visualizations to explore heterogeneous tax data, facilitating the investigation of suspicious transactions among related taxpayers.

At the national level, the significant work of da Silva et al [2] deserves attention for its exploration of Bayesian networks in the context of the Brazilian federal Tax Administration. This study focuses on the application of Bayesian networks for the targeted scrutiny of Income Tax declarations, enhancing the efficiency of the tax audit process. Additionally, Xavier et al [8] contributes to the field by introducing an innovative methodology for identifying potential tax evaders. Their approach leverages open and public data, employing Random Forest, Neural Networks, and Graphs to construct a binary classification model. This model effectively distinguishes between "default" and "reputable" company profiles specifically within the State of Goiás.

### 3 PROPOSED SOLUTION

To address the presented challenge and notably contribute to the initial phase of the selection process, this study proposes the application of supervised learning algorithms to initially predict the annual income of an individual taxpayer (dependent variable – Y) utilizing various internal data provided to the Brazilian Tax Authority (RFB) in the taxpayer's own or third-party declarations (independent variables – X).

Conversely, the indication of omitted income would be characterized by the positive difference between the annual income value estimated by the predictive model and the value actually declared in the taxpayer's Individual Income Tax Return (DIRPF). This disparity will be henceforth referred to as **Estimated Omission Income (EOI)**.

$$EOI = \max(0, \text{predicted} - \text{reported value})$$

The method focus its analysis on the taxpayers whose EOI surpasses the **MAE (Mean Absolute Error)** metric of the predictive model. Insignificant deviations are deemed inconsequential, as they fall within an acceptable range of potential prediction error for the model and/or do not manifest a substantial degree of omission warranting the cost associated with a tax audit.

In consideration of this, taxpayers are stratified into three delineated groups based on the computed EOI:

- **Group 1** (predicted value > declared value): Instances where the discrepancy between the predicted and declared values exceeds the threshold of the MAE. This is our focal group.
- **Group 2** (predicted value  $\approx$  declared value): variances between predicted and declared values within the specified range (-MAE, +MAE), which characterizes the behavioral pattern of the "compliant taxpayer."
- **Group 3** (predicted value < declared value): This category encapsulates scenarios wherein taxpayers declare an annual income exceeding the predicted value.

In order to further fine-tune the selection process, aligning it with the constrained workload capacity and mitigating the incidence of false positives, the concluding step of the proposed method entails the implementation of statistical techniques for outlier detection within the identified focal group (Group 1). This strategic approach is designed to prioritize an in-depth individual analysis of these taxpayers.

In essence, the method comprises the following steps:

- Step 1** Employing supervised learning to predict the annual income of an individual taxpayer.
- Step 2** Computing the **EOI** by deriving the difference between the predicted and actual (reported) values.
- Step 3** Selecting the group of taxpayers for whom  $EOI > MAE$ .
- Step 4** Identifying outliers within this group and arranging them in descending order based on **EOI**.

The model can be regarded as a "selection rule" whose objective is to generate a preliminary list of taxpayers ordered in descending order by the estimated omission income value (EOI). This list serves as a prioritized queue for the individual analysis of taxpayers in the subsequent stage of the audit case selection process.

### 3.1 Methodology

This study was fundamentally based on the CRISP-DM framework [7], a widely acknowledged and comprehensive approach to data mining. CRISP-DM comprises six phases – business understanding, data understanding, data preparation, modeling, evaluation and deployment. It is an iterative process, allowing for revisiting earlier phases as new insights emerge, enabling continuous enhancement and fine-tuning of the analysis.

## 4 EXPERIMENT

### 4.1 Business Understanding

The RFB's Individual Taxpayer Registry contains information on over 200 million Individual Taxpayer Identification Numbers (Cadastro da Pessoa Física – CPF). Despite widespread compliance, the sheer volume of absolute numbers highlights the challenge of navigating this vast dataset and prioritizing cases with economic potential in the selection process.

In addressing this challenge, the application of data mining techniques emerges as a valuable alternative, providing a sophisticated means to sift through the vast dataset and identify cases warranting closer scrutiny, considering both their economic impact and the likelihood of tax evasion. For the purpose of establishing a proof of concept, the scope of this study was delimited as follows:

- **Type of Taxpayer:** Individual

- **Calendar Year (CY):** 2019
- **Occupational Group:** Member or public servant of the federal direct administration

The scoping limitation was deemed necessary to reduce the size of the dataset, originally comprising tens of millions, to approximately **400.000 rows**. This reduction facilitates a more in-depth analysis of data pertaining to taxpayers within the same group, characterized by greater homogeneity. Moreover, a deliberate choice was made to focus exclusively on datasets associated with taxpayers whose disclosed occupational nature in the DIRPF is classified as "*21 - Member or public servant of the federal direct administration*". This particular subgroup furnishes annual income values reported by a dependable third party, specifically, the Federal Government. It is presumed that this subset inherently harbors more trustworthy information regarding the annual income (target variable), rendering it particularly well-suited for utilization as input to the machine learning model.<sup>1</sup>

## 4.2 Data Understanding

The repository of taxpayer data accessible to the RFB is vast, encompassing data on the scale of thousands of terabytes distributed across several databases consolidated within a comprehensive data-lake infrastructure. Since the primary objective of the model is to predict the annual income of a taxpayer, this study prioritized those variables that could potentially contain information regarding the socio-economic status of the taxpayer.<sup>2</sup>

The results of queries across the databases of interest were consolidated, using the CPF and calendar year (CY) as merging criteria. The outcome was subsequently exported to a unified file, comprising **395.560 rows** (records) and **59 columns** (features). Each row contains information related to a unique CPF. The pertinent features for this study are described in Table 1<sup>3</sup>.

The first four variables are of text type, while the remaining ones are numeric. It is worth clarifying that the selection of the aforementioned variables was made because they represent, to a greater or lesser extent, indicators of income and/or wealth of a taxpayer. It is expected, therefore, that they have some relationship with the target variable intended to be estimated (*vl\_income*).

## 4.3 Data Preparation

The following procedures and techniques were applied to convert the data into a format suitable for modeling and enriching the model with more relevant information.

**4.3.1 Missing Values.** Missing values in several columns, such as *qt\_emp*, *qt\_resp*, *vl\_buy*, and *vl\_cred*, were replaced with 0 to signify the absence of specific events. This approach aligns with domain knowledge and appropriately captures the absence of occurrences.

<sup>1</sup>The detailed information about the compensation of federal public servants is openly accessible to the public via the Federal Government Transparency Portal – <https://portal.datatransparencia.gov.br/servidores>

<sup>2</sup>In compliance with legal confidentiality requirements for tax information, this study used aggregated and/or anonymized data, ensuring the prevention of direct or indirect identification of individual taxpayers. This deidentification process does not impede the comprehension, development or reproducibility of the study

<sup>3</sup>DIRF – Income Tax Withholding Statement reported by the withholding agent. SPED-NFe – Digital Accounting System for Electronic Invoices. DECRED – Statement of credit card transactions reported by credit card issuers.

**Table 1: Description of the main features**

Feature	Description
cd_ocu	Primary occupation as reported in DIRPF.
age	Age of the taxpayer
gender	Gender of the taxpayer
mar_stat	Marital status
type	Type of taxpayer according to the Tax Administration internal classification
qt_dep	Number of dependents reported in DIRPF
qt_comp	Number of companies in which the taxpayer is registered as a shareholder in the QSA (Shareholders and Administrators Registry).
qt_resp	Number of companies where the taxpayer is listed as a responsible party in the QSA.
vl_income	Annual income amount reported in DIRPF (target variable).
vl_asset	Value of assets and rights reported in DIRPF
vl_dirf	Annual income amount reported in DIRF by the withholding entity
vl_fin	Total value of credits in the account as reported by the financial institution(s).
vl_buy	Total value of electronic purchase invoices reported in SPED NF-e.
vl_cred	Value reported in DECRED by the credit card operator, referring to the total credit card transactions carried out by the in a specific year.

However, for columns *vl\_dirf* and *vl\_cred*, which exhibited a small proportion of missing values (less than 0.01%), the strategy involved removing entire rows with missing data. This decision was also based on the random and unpatterned distribution of the missing values across the dataset.

**4.3.2 Data Transformation.** As part of the data transformation phase dedicated to preparing the dataset for the machine learning model, we encoded categorical variables by converting them into numerical representations using the one-hot encoding technique. Additionally, we applied the normalization technique to numeric variables, scaling them to a standardized range between 0 and 1. Moreover, we implemented discretization to transform continuous variables into discrete or categorical intervals. An illustrative example of this was the transformation of the *age* variable into an age range.

**4.3.3 Feature Engineering.** Feature engineering involves creating new features, modifying existing ones, and selecting the most relevant variables to improve the model's ability to make accurate predictions. In pursuit of heightened predictive capabilities, the features in Table 2 have been derived from the original set, strategically designed to confer additional predictive power to the machine learning model.

**4.3.4 Outliers.** To further tailor the data to enhance the outcomes of the ML model, outlier values were removed. For this purpose, values with a z-score exceeding 6 (six) standard deviations from the



**Table 2: New features generated during Feature Engineering**

Feature	Description
vl_var	vl_asset (2019) - vl_asset (2018).
qt_qsa	qt_comp + qt_resp
rat_1	Ratio between vl_fin e vl_dirf
rat_2	Ratio between vl_var e vl_dirf

population mean were deemed outliers. A relatively high z-score threshold was chosen to retain as much information as possible from extreme values. It is worth noting that the removal of outliers from the dataset proved essential in reducing the variance of test scores during the application of the k-fold cross-validation technique for the machine learning model.

After the data preparation stage, the final number of records was reduced to **390.165**. A correlation heatmap, as shown in Figure 1 was generated to visually represent the Pearson correlation coefficients among the various features in the dataset.

**Figure 1: Pearson Correlation Heatmap**

## 4.4 Modeling

To tackle the presented challenge, we embraced the supervised machine learning paradigm, specifically employing regression methods to predict the dependent variable (*vl\_income*) based on independent variables. The selected algorithm was Gradient Boosting, implemented by the *xgboost*[1] Python package, strategically designed for enhanced speed and performance.

**4.4.1 Hyperparameter tuning.** The hyperparameter tuning process employed the *hyperopt*<sup>4</sup> Python package, implementing Bayesian methods for sophisticated optimization. After extensive experimentation with different hyperparameter values, default settings were retained, as substantial changes in model accuracy were not observed with alternative configurations.

**4.4.2 Metrics.** For a comprehensive assessment of the model's performance, in this study, we chose to employ two different metrics: R2 score and MAE.

The **R2 score** is a common metric for assessing regression models, indicating the proportion of variability in the dependent variable predicted by the model based on independent variables. Ranging

<sup>4</sup><https://hyperopt.github.io/hyperopt/>

from 0 to 1, a score of 1 signifies a perfect fit, explaining all variations, while 0 suggests a failure to explain any variation.

The **MAE** metric calculates the average of the absolute differences between the model predictions and the actual values. It was chosen for its ease of interpretation and allows for a direct comparison with the target variable (*vl\_income*), as both are in the same unit of measurement.

**4.4.3 Baseline.** To evaluate the model, in addition to the metrics, we found it important to establish a baseline parameter for comparison. The baseline is nothing more than a naive guess against which we can compare the model results. For this case, we consider a reasonable baseline to be the median value of the target variable we aim to predict. Applying this concept, if we adopt the median guess (178.310,38) to predict the income value for all instances in the dataset, the Mean Absolute Error would be 116.220,36. It is expected that the supervised model should significantly reduce this error.

Baseline parameter (median) :	<b>178.310,38</b>
MAE performance on the baseline:	<b>116.220,36</b>

## 5 RESULTS

The dataset was split into 80% for the training set and 20% for the test set. Table 3 summarizes the results achieved by the proposed model in the selected metrics.<sup>5</sup>

**Table 3: Performance metrics results**

Metric	Value
MAE	23.356,21
score_train	0.90
kfold_mean	0.88
kfold_std	0.01

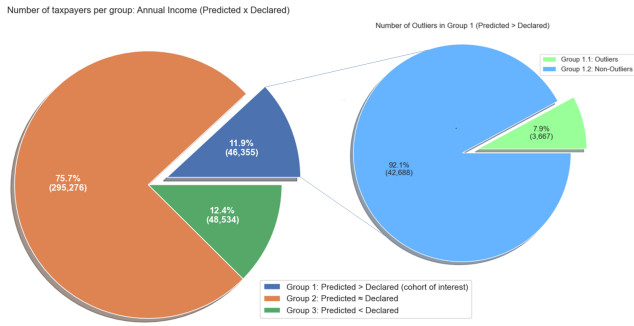
The model demonstrated a remarkable forecasting ability, with an average coefficient of determination (R2) close to 0.88, obtained through the *k-fold* cross-validation method with k=5. This value suggests a considerably high accuracy in the model's estimates. Additionally, it is relevant to mention the low standard deviation value which indicates that the model's performance is highly stable and consistent, regardless of variations in the datasets used.

Furthermore, considering the MAE metric, the model achieved a noteworthy value of **R\$ 23.356,21**. This result is significantly lower than the MAE of the reference model (baseline), which stood at **R\$ 116.220,36**. This contrast underscores the model's promising predictive ability, exhibiting an approximately fivefold improvement in precision compared to the baseline. This substantial difference clearly illustrates the positive impact of the model on the forecasting process, highlighting its efficacy in producing more accurate and reliable estimates.

The graph in Figure 2 illustrates how taxpayers were classified by EOI according to the criteria described in Section 3.

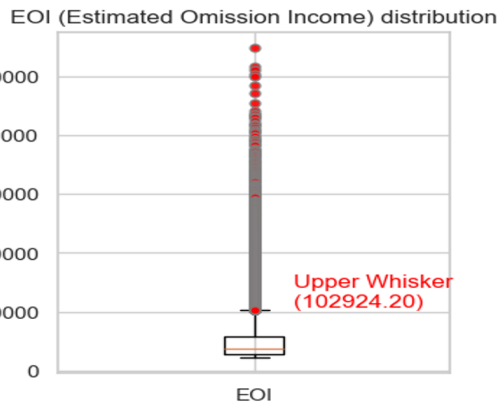
The identified Group 1, displaying signs of income omission according to our model's classification, constitutes only **11.9%** of the

<sup>5</sup>*score\_train*: R2 score on the training set. *kfold\_mean* and *kfold\_std* are the mean/standard deviation of R2 score obtained by applying the k-fold method (k=5).



**Figure 2: Number of taxpayers per group and number of outliers in Group 1**

initial taxpayer population. Although this initial filter significantly reduces the dataset, there are still **46,355** taxpayers necessitating further analysis in the subsequent individual assessment stage. In essence, further refinement of the list is required. In this context, the next step in the proposed method is to identify outliers within the group of interest (Group 1). To assist in this task, we employed the **boxplot** method, which is a type of chart that displays the distribution of numerical data and skewness by presenting the five-number summary of a dataset: including the minimum score, first (lower) quartile, median, third (upper) quartile, and maximum score. Box plots are valuable as they reveal outliers within a dataset, defined as data points located outside the whiskers of the box plot. In this case, our focus is on points above the upper whisker, calculated as  $Q3 + 1.5 * IQR$  (interquartile range). Examining the boxplot (Figure 3) corresponding to Group 1 (predicted > declared), we observe several points above the upper whisker, with a value of **R\$ 102,924,20**, which was considered as the threshold to define outliers.

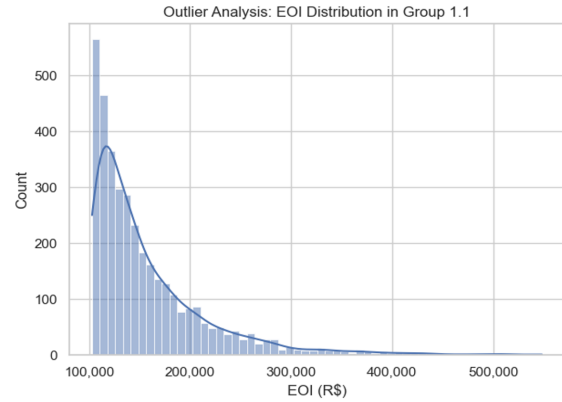


**Figure 3: Boxplot of EOI Distribution**

Given the outlier classification, Group 1 was further divided into two subgroups, as shown in Figure 2 :

- **Group 1.1 (outliers):** Taxpayers with EOI > 102,924,20
- **Group 1.2 (non-outliers):** Taxpayers with EOI ≤ 102,924,20

The proposed framework was able to refine the initial taxpayer list by applying an objective and impartial "filter", leading to the pre-selection of **3,667** taxpayers (Group 1.1), constituting around **0.95%** of the initial universe of taxpayers under analysis. The histogram in Figure 4 details the frequency of EOI values in **Group 1.1**.



**Figure 4: Outlier Analysis: EOI Distribution in Group 1.1**

Upon closer examination, we observe that, among the taxpayers in Group 1.1 (outliers), the highest EOI is **R\$ 548,391,88** and the total sum of EOI reaches **R\$ 572,791,819,89**. In a hypothetical scenario where the model predictions are entirely accurate, we would be looking at an estimated income omission value of approximately half a billion Brazilian reais (BRL). This is a remarkably substantial amount, especially considering that we are dealing exclusively with individual taxpayers whose occupation falls under the category of federal public servants, representing a small fraction of the overall taxpayer population.

The ultimate result of the proposed audit case selection process is a descending ordered list based on EOI values from taxpayers in Group 1.1, which represents the cohort with the highest risk of tax evasion and, hence, should be prioritized for further scrutiny.

## 5.1 Explainability

*Explainable Artificial Intelligence (XAI)* strives to enhance the transparency and interpretability of AI models. In the realm of selecting taxpayers for audit, XAI plays a pivotal role in ensuring fairness, compliance, and trust in machine learning models. It addresses legal and ethical considerations, fostering responsible AI use within tax administration. XAI facilitates a deeper understanding and justification of ML decisions, bolstering confidence and acceptance of AI technologies in tax-related processes.

For this study, we opted for *SHAP (SHapley Additive exPlanations)*[5], a prominent XAI framework rooted in cooperative game theory. SHAP provides a unified measure of feature importance for each prediction, contributing to a clearer interpretation of machine learning outputs. The SHAP summary plot presented in Figure 6 illustrates the influence of each feature on the model's output across the entire dataset. This visual aid assists users in comprehending the overall behavior of the model, shedding light on how individual features contribute to the model predictions. While understanding

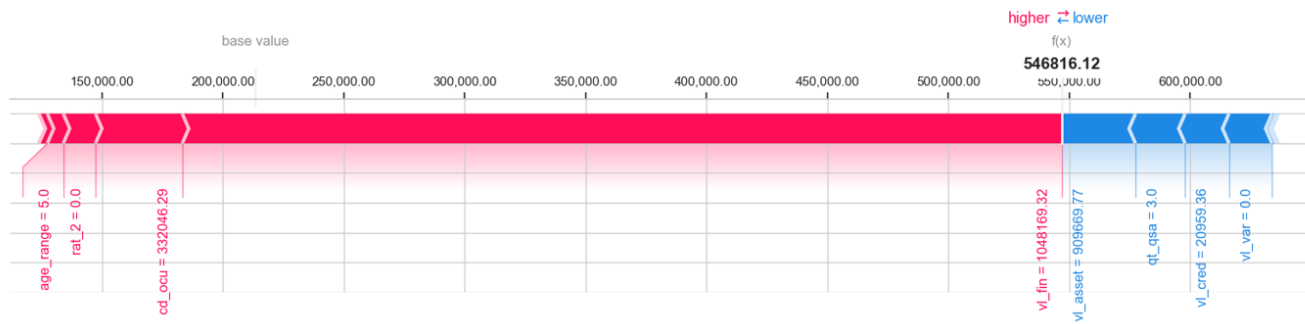


Figure 5: Local Explainability – SHAP values for the particular instance with the highest EOI.

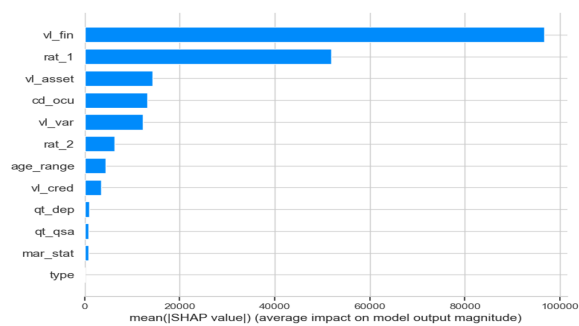


Figure 6: Global Explainability – Feature Importance Based on Mean SHAP Values

the global importance of features is crucial, comprehending why a particular taxpayer is selected for audit is equally vital. The SHAP framework also provides local explainability, revealing specific factors influencing the model's decision for a particular prediction. Figure 5 showcases the SHAP values for the highest Estimated Omission Income (EOI) instance, providing insight into the model's prediction for this specific taxpayer.

The SHAP plot visually emphasizes the features that had the greatest impact on this outcome. Positive SHAP values, represented in red, indicate features that increased the prediction, while negative values, displayed in blue, suggest features that decreased the prediction. These features may signify different levels of atypical behavior, potentially serving as indicators of risk factors.

## 6 CONCLUSION

In this study, we have proposed a novel approach to enhance the audit case selection process in the Brazilian Tax Administration. We have applied supervised learning and data mining techniques to a real-world dataset reflecting taxpayers' socioeconomic situations. The application of advanced Machine Learning, notably explainable AI, not only propels data analysis in tax administration but also paves the way for future research and audit improvements while addressing legal and ethical considerations.

The outcomes of this work have the potential to revolutionize the RFB's taxpayer selection process by integrating Artificial Intelligence methods into the traditional rule-based system leveraging

the use of big data. Beyond practical impacts, the research enriches academic literature on AI's role in combating tax fraud and evasion, filling a crucial gap due to the complexity of working with real-world data. Its diverse contributions extend to influencing audit effectiveness, promoting justice and the overall integrity of the fiscal system.

Potential avenues for future study expansion include: conducting individual analyses on a subset of Group 1.1 instances to confirm or refute indications of income omission and validate the model; testing the model on taxpayers beyond the original scope, particularly those with occupations other than "federal public servant"; incorporating new independent variables available in the RFB's internal databases, not considered in this study, to improve the model's predictive capacity; assessing the creation of derived features based on existing variables; and employing unsupervised outlier detection algorithms to compare and analyze the results obtained.

## REFERENCES

- [1] Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (San Francisco, California, USA) (KDD '16). ACM, New York, NY, USA, 785–794. <https://doi.org/10.1145/2939672.2939785>
- [2] da Silva L. S., Rigitano H. de C., Carvalho R. N., and Souza J. C. F. 2016. Bayesian networks on income tax audit selection — a case study of Brazilian tax administration. In *Bayesian Modeling Application Workshop (BMAW)*.
- [3] D. de Roux, B. Perez, A. Moreno, M. D. P. Villamil, and F. Figueroa. 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 (KDD '18)*. 215–222. <https://doi.org/10.1145/3219819.3219878>
- [4] Y. Lin et al. 2021. TaxThemis: Interactive Mining and Exploration of Suspicious Tax Evasion Groups. *IEEE Transactions on Visualization & Computer Graphics* 27, 02 (2021), 849–859. <https://doi.org/abs/2009.03179>
- [5] Scott M Lundberg and Su-In Lee. 2017. A Unified Approach to Interpreting Model Predictions. In *Advances in Neural Information Processing Systems* 30, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.). Curran Associates, Inc., 4765–4774. <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>
- [6] Miloš Savić et al. 2021. Tax Evasion Risk Management Using a Hybrid Unsupervised Outlier Detection Method. <https://arxiv.org/pdf/2103.01033.pdf>
- [7] Rudiger Wirth and Jochen Hipp. 2000. Crisp-dm: Towards a standard process model for data mining. (2000), 29–39.
- [8] Otávio Xavier et al. 2022. Tax evasion identification using open data and artificial intelligence. *Revista de Administração Pública* 56 (2022), 426–440. <https://doi.org/10.1590/0034-761220210256x>
- [9] M. Zumaya et al. 2021. *Identifying Tax Evasion in Mexico with Tools from Network Science and Machine Learning*. Springer. [https://doi.org/10.1007/978-3-030-81484-7\\_6](https://doi.org/10.1007/978-3-030-81484-7_6)