# Tax Prediction System using machine learning

**Sanket Lade  Shubham Shinde  Yuvraj Sutar  Chirag Ramraje**

**Guide name: Prof. V. D. Badgujar**

**Department of Computer Engineering Brahma valley college
of engineering & research institute
Nashik**

-----------------------------------------------------------------***-----------------------------------------------------------------

**Abstract**        Machine learning gives strategies, tools and equipment, which assist to learn mechanically and to make correct predictions based totally on beyond observations. The records are retrieved from the actual time environmental setup. Machine getting to know techniques can help in the integration of laptop-primarily based structures in predicting the dataset and to improve the performance of the device. The main reason of this paper is to provide Tax Predictions from given data and also fraud detection. Such contrast helps to offer the correct result in algorithms.

For this reason evaluating, it tries to determine tax benefits which are more likely to be utilized by ability fraud taxpayers by means of investigating the non-public income tax structure. Secondly, it targets at characterizing thru socioeconomic variables the phase profiles of potential fraud taxpayer to offer an audit selection approach for enhancing tax compliance and improve tax design. Random forest algorithms is a tedious undertaking, for real time dataset. The combination of statistics Feature Extraction proposed gives precious statistics to contribute to the examine of tax fraud.

*Key Words***:** :-  Machine-learning Techniques, Random Forest Algorithm, Audit Selection        Strategy, Tax audit, Feature Extraction.

**INTRODUCTION  :**

Earnings tax is an important source of revenue to government in both growing and evolved nations. The amount of revenue to be generated by government from such taxes for its expenditure programmers relies upon, among different things, at the willingness of the taxpayers to conform with the tax legal guidelines of a rustic. Taxation is a way through which governments finance their expenditure by implementing expenses on citizens and corporate entities. There are one of a kind styles of tax, however simplest the important one, that this take a look at focuses on, specifically countrywide tax (non-public profits tax).Whilst term analytics is often utilized by tax practitioners, it's far a wide time period, used to describe the entirety from business intelligence, dashboards, predictive and prescriptive tax analytics, to extra superior areas including system learning (ml), information mining.

Device gaining knowledge of offers strategies strategies and equipment, which help to study routinely and to make accurate predictions based totally on beyond observations. Device studying is popularly being utilized in areas of commercial enterprise like statistics analysis, financial evaluation; stock market forecast and so on. Classification is used to build category tree for predicting non-stop established variables and specific predictor variables. Tax fraud detection entails processing a big quantity of facts searching for fraudulent

behavior that calls for speedy and green algorithms, among which facts mining presents relevant strategies that can help tax administration to take preventive measures and improve tax design. Auditing tax declarations is a gradual and luxurious procedure, in order that, tax government required to broaden fee-efficient techniques to tackle this hassle and improve tax layout. This trouble motivates our thought. In our analysis we explore the applicability of the records mining strategies in developing a segmentation version which can make contributions to tax design evaluation and the characterization of the segments of capability fraud taxpayers inside the non-public income tax. In spite of the increase within the use of these screening and type models for detecting fraud styles orientated at audit making plans, there are no studies that target the identification of tax blessings within the earnings tax structure which are more likely to be used by ability fraud taxpayers.

**AIM OF PROJECT**:

   Maximum of methods are mainly supervised mastering, or rely upon the past conduct of taxpayers. On this instance, they use marked statistics indicating a fraud and use these statistics to create each prediction and classifications models. Consequences of those strategies are best, however those strategies can´t be normally applied across tax fraud given that statistics is not without problems available. This take a look at objectives to fill this hole with the aid of growing a popular method primarily based on a aggregate of foremost components analysis, neural networks and selection timber that lets in for the detection of over-claiming tax benefits and which can be implemented to extraordinary forms of taxes without the want to have get entry to to tax fraud categorised ancient facts. Its intention is to help tax audit professionals in defining essential factors to take into account when appearing detection of fraudulent taxpayers.

We here show that the proposed machine outperforms present statistical methods to tax default prediction

 By measuring the prediction electricity of the economic indicators, this look at additionally examines their importance and establishes a comprehensive early-caution gadget concerning the status of company tax payment.

**Need of Project**  To cope with the economic significance of unpaid taxes by means of the use of an automated gadget for predicting a tax default. Too little attention has been paid to tax default prediction inside the beyond. World bank information claim that about 40% of companies around the globe pay their taxes however 60% fail to pay their taxes, and those amounts may not be recovered at some point of upcoming tax years. The information additionally record that prices of tax defaults are increasing global. Considering the financial importance of unpaid taxes. Little research has been carried out in predicting the tax status of firms.

        Monetary fraud is an important difficulty that incurs fees in terms of the loss of presidency revenues, which leads to less green tax applications and the inequity among evaders and sincere filers. Tax management are underneath growing pressure, because the monetary crisis of 2008 and the massive deficits that accompanied, to collect extra tax sales and decrease monetary fraud. Effective control of tax fraud requires addressing a fundamental statistical problem of non-detection, which could bias estimates of the overall amount of fraud and the relative fraud propensities of different socioeconomic organizations.

# Software Requirement

➢ **Python 3.6:**

   **Packages**

     Flask
     Open cv
     NLTK Tesnsor flow & Keras

IDLE is Python's Integrated Development and Learning Environment. IDLE has the following features:
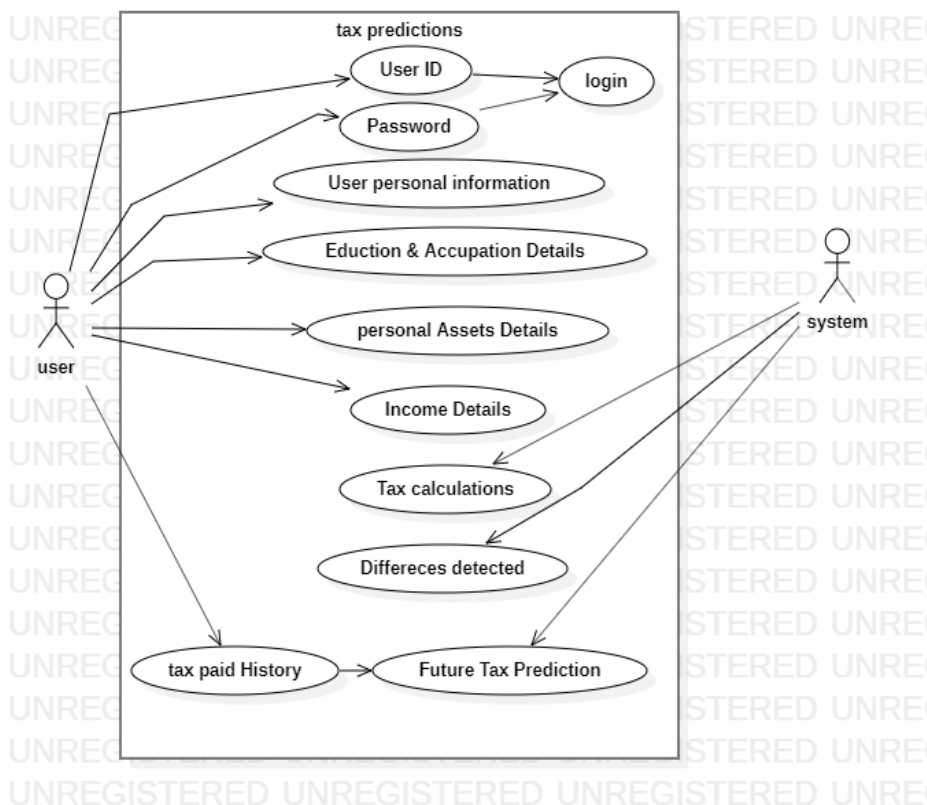
- coded in 100% pure Python, using the tkinter GUI toolkit
- cross-platform: works mostly the same on Windows, Unix, and macOS
- Python shell window (interactive interpreter) with colorizing of code input, output, and error messages
- multi-window text editor with multiple undo, Python colorizing, smart indent, call tips, auto completion, and other features
- search within any window, replace within editor windows, and search through multiple files (grep)
- debugger with persistent breakpoints, stepping, and viewing of global and local namespaces.
-
  **SYSTEM ARCHITECTURE:**

## Existing System :

In the existing system for taxation analysis using classification techniques the data mining algorithm used in machine learning. Data mining has many existing and potential applications in tax administration. The data under analysis represent the income and taxation particulars. The income data pertaining to the financial year ending 31st march, 2006 is segregated under various Heads and Gross total Income, Deductions and Net Taxable Income along with Income tax payable thereon and interest, if any, are also listed. Under Income tax Act, 1961, Persons or entities, known as Assesses, earning Income are divided in to various categories.

**Use Case Diagram**. A use case diagram at its simplest is a representation of a user's interaction with the system and depicting the specifications ract with the system. User upload image on browser. System compare this image with previous dataset and produce result of a use case. A use case diagram can portray the different types of users of a system and the various .
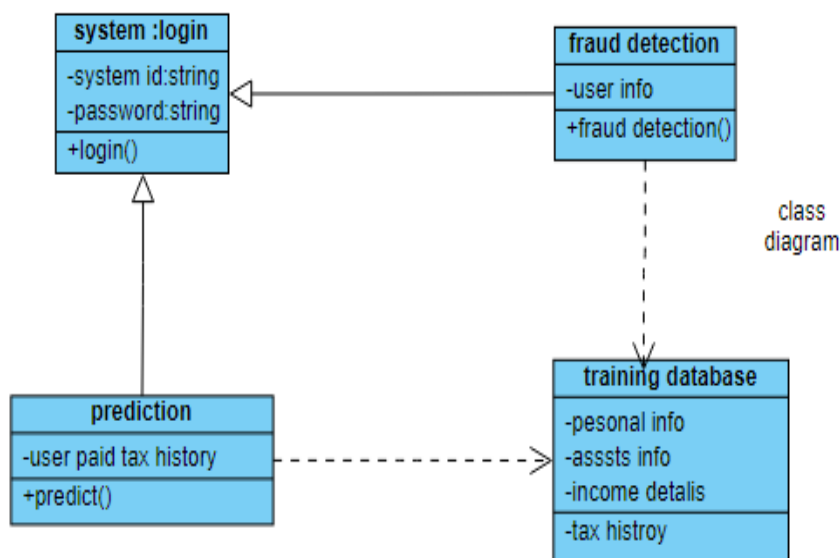
**Data Flow Diagram**

DFD describes the processes that are involved in a system to transfer data from the input to the file storage and reports generation. Data flow diagrams can be divided into logical and physical. The logical data flow diagram describes flow of data through a system to perform certain functionality of a business. The physical data flow diagram describes the implementation of the logical data flow.

**Class Diagram**

The class diagram is the main building block of object-oriented modelling. It is used for general conceptual modelling of the structure of the application, and for detailed modelling translating the models into programming code. Class diagrams can also be used for data modelling. The classes in a class diagram represent both the main elements, interactions in the application, and the classes to be programmed.
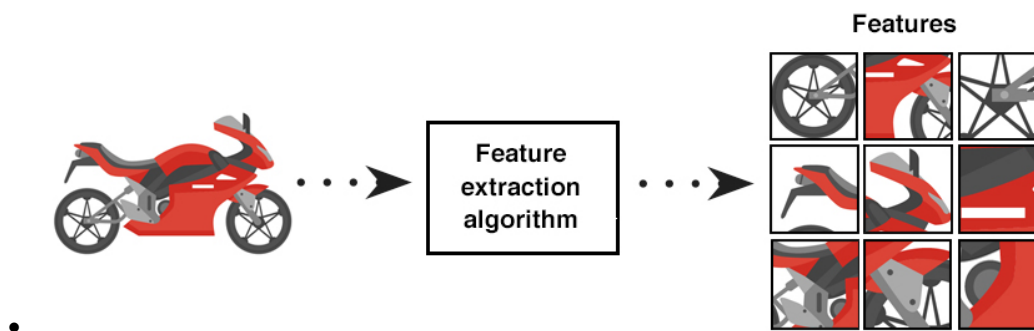


CLASS DIGRAM

**Feature extraction:**

Feature extraction is a core component of the computer vision pipeline. In fact, the entire deep learning model works around the idea of extracting **useful features** which clearly define the objects in the image. We're going to spend a little more time here because it's important that you understand what a feature is, what a vector of features is, and why we extract features. A feature in machine learning is an individual measurable property or characteristic of a phenomenon being observed. Features are the input that you feed to your machine learning model to output a prediction or classification. Suppose you want to predict the price of a house, your input features (properties) might include: square foot, number_of_rooms, bathrooms, etc. and the model will output the predicted price based on the values of your features. Selecting good features that clearly distinguish your objects increases the predictive power of machine learning algorithms.

**Working of Random Forest Algorithm**

We can understand the working of Random Forest algorithm with the help of following steps −

- **Step 1** − First, start with the selection of random samples from a given dataset.

- **Step 2** − Next, this algorithm will construct a decision tree for every sample. Then it will get the prediction result from every decision tree.

- **Step 3** − In this step, voting will be performed for every predicted result.

- **Step 4** − At last, select the most voted prediction result as the final prediction result.

-



**FEATURE EXTRACTION ALGORITHM**

In machine learning projects, we want to transform the raw data (image) into a *features vector* to show our learning algorithm how to learn the characteristics of the object.In the image above, we feed the raw input image of a motorcycle to a feature extraction algorithm. Let's treat the feature extraction algorithm as a black box for now and we'll come back to it soon. For now, we need to know that the extraction algorithm produces a *vector* that contains a list of features. This is called **features vector** which is a 1D array that makes a robust representation of the object. It is important to call out that the image above reflects features extracted from just one motorcycle. A very important characteristic of a feature is **repeatability**.

**COCLUSION**

In this project we successfully implemented fraud detection system in tax. The results obtained in this study present a wide range of possibilities to the improve tax fraud detection, through the use of the kind of predictive tools to find fraud patterns which could be described a priori, through sensitivity analysis. Also we predicted how much future tax should be paid by person using feature extraction and random forest algorithm.

➢ **Future Scope**

➢ As a future work it is planned to apply various Machine Learning and IOT approaches to large-scale real data provided by tax authorities.
➢ In the future, it would be of great interest to realize applications of this methodology in other taxes.

## REFERENCES:

➢

➢ [1] Abdallah, A.; Mohd, A.M.; Anazida, Z. Fraud detection system: A survey. J. Netw. Comput. Appl. 2016, 68, 99–113. [CrossRef]

➢ [2] Anyaeche, C.O.; Ighravwe, D.E. Predicting performance measures using linear regression and neural network: A comparison. Afr. J. Eng. Res. 2013, 1, 84–89.

➢ [3] Y. Jiang and S. Jones, ''Corporate distress prediction in China: A machine learning approach,'' Accounting Finance, vol. 58, no. 4, pp. 1063–1109, Dec. 2018.

➢ [4] J. Vanhoeyveld, D. Martens, and B. Peeters, ''Value-added tax fraud detection with scalable anomaly detection techniques,'' Appl. Soft Comput., vol. 86, Jan. 2020, Art. no. 105895.

➢ [5] P. Hajek and K. Michalak, ''Feature selection in corporate credit rating prediction,'' Knowl.-Based Syst., vol. 51, pp. 72–84, Oct. 2013.

➢ [6] W.-C. Lin, Y.-H. Lu, and C.-F. Tsai, ''Feature selection in single and ensemble learning-based bankruptcy prediction models,'' Expert Syst., vol. 36, no. 1, Feb. 2019, Art. no. e12335.

➢ [7] H. Höglund, ''Tax payment default prediction using genetic algorithmbased variable selection,'' Expert Syst. Appl., vol. 88, pp. 368–375, Dec. 2017.

➢ [8] J. Ruan, Z. Yan, B. Dong, Q. Zheng, and B. Qian, ''Identifying suspicious groups of affiliated-transaction-based tax evasion in big data,'' Inf. Sci., vol. 477, pp. 508–532, Mar. 2019.

➢ [9] W. Didimo, L. Grilli, G. Liotta, L. Menconi, F. Montecchiani, and D. Pagliuca, ''Combining network visualization and data mining for tax risk assessment,'' IEEE Access, vol. 8, pp. 16073–16086, 2020.

➢ [10] D. Marghescu, M. Kallio, and B. Back, ''Using financial ratios to select companies for tax auditing: A preliminary study,'' in Organizational, Business, and Technological Aspects of the Knowledge Society. Berlin, Germany: Springer, 2010, pp. 393–398.

➢ [11] M. Moscatelli, F. Parlapiano, S. Narizzano, and G. Viggiano, ''Corporate default forecasting with machine learning,'' Expert Syst. Appl., vol. 161, Dec. 2020, Art. no. 113567.