

# Forecasting Cybercrime Using Machine Learning Algorithms in Python

Durga.C<sup>1</sup>, Dr.Sreejith Vignesh B P<sup>2</sup>

<sup>1</sup>Junior Researcher Department of Information Technology, Sri Krishna Adithya college of arts and science

<sup>2</sup>Associate Professor & Head, Department of Information Technology, Sri Krishna Adithya college of arts and science

**Abstract-** Cybercrime has increased dramatically as a result of the fast development of digital technologies, which represents a serious risk to global security. Traditional reactive methods frequently fall behind the changing nature of cyberattacks. Utilizing Python within the Google Colab cloud environment and implementing cutting-edge Machine Learning (ML) algorithms, this study seeks to predict cybercrime trends. Using a dataset of local cyber events, the research uses Multi-Output Regressors and Random Forest Classifiers to discover underlying patterns and forecast particular future threats. The study shifts from simple detection to predictive analytics, providing a scalable framework for businesses and law enforcement. The results show that integrating web-based frontends (HTML/CSS/JS) with ensemble-based learning produces an accessible and extremely precise approach to proactively reduce risk.

**Keywords:** Cybercrime, Forecasting, Machine Learning, Python, Google Colab, Predictive Modeling, Random Forest.

## I.INTRODUCTION

The dependence on technology and the internet has grown significantly in all industries in today's digital age, including banking, healthcare, education, and e-commerce. Although this digital revolution has increased convenience and efficiency, it has also resulted in a large increase in cybercrime. Cybercrimes, which range from ransomware assaults, identity theft, and phishing to data breaches on a massive scale, put individuals, governments, and multinational corporations at risk. The intricacy and It becomes more and more difficult to identify and prevent these crimes due to their increasing complexity. [2]

Most conventional cybersecurity systems are reactive,

concentrating on reacting to assaults after they have taken place. Nonetheless, given the constantly changing methods used by cybercriminals, there is a pressing need for proactive measures that can foresee and mitigate potential threats before they cause harm.[3] Given the context, machine learning (ML) has proven to be a potent instrument for analyzing massive amounts of cyber-related data information, identifying underlying trends, and predicting upcoming cyber risks with great precision.

The goal of this study article is to examine how machine learning techniques may be used in the Eclipse Integrated Development Environment (IDE) to predict cybercrime. [4] Eclipse offers a versatile and robust platform for integrating ML libraries and frameworks, facilitating the efficient training, testing, and deployment of models.[5] The study uses algorithms like Decision Trees, Random Forests, Support Vector Machines, and Neural Networks to try to determine how well they can forecast cybercrime trends.

The main goals of this study are to illustrate how predictive analytics can, compare the performance of different machine learning models in predicting cyber threats, and offer insights into developing more proactive and intelligent defense systems. It is anticipated that the findings of this research will help legislators, law enforcement officers, and cybersecurity experts in identifying possible threats and minimizing the effects of cybercrime in the future.[7]

## II.PROBLEM MOTIVATION WITH REAL - WORLD STATISTICS

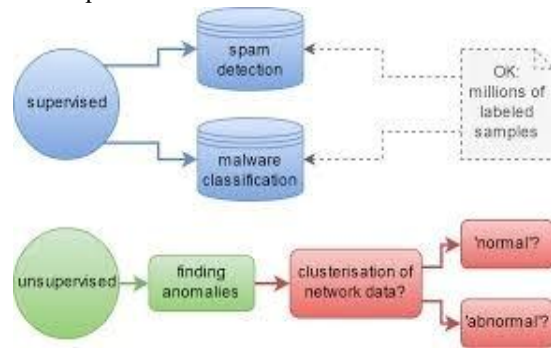
One of the biggest risks to the worldwide digital

economy is cybercrime. According to Cybersecurity Ventures reports, global cybercrime damage expenses are predicted to amount to \$10.5 trillion per year by 2025, making it one of the biggest economic offenses ever [8]. The rise is just as frightening in the Indian setting; according to complaints made to the National Cyber Crime Reporting Portal, there was an astounding 200% increase between 2019 and 2023 [10]. This data shows a large increase in financial fraud, identity theft, and online scams, which varies greatly by geographic area and administrative area. The impetus behind this initiative is that current security measures are primarily reactive, addressing a breach only after the harm has been done [12]. This delay causes significant operational downtime and financial losses. Additionally, several earlier studies prioritize worldwide trends over regional trends. By progressing toward a proactive forecasting approach, our study fills this gap. To forecast future patterns, we employ Python-based Multi-Output Regression models in a Google Colab environment to analyze historical district-wise data (specifically the 2022 datasets). In contrast to a binary "attack or no attack" prediction, this project predicts the severity and kind of crime likely to take place in a particular area, such as phishing, financial fraud, or identity theft. Using Machine Learning to analyze these historical trends allows law enforcement and legislators to allocate resources efficiently, thereby improving regional cyber defense's efficiency, adaptability, and resilience [15].

### III.LITERATURE REVIEW & REVIEW OF RECENT RELATED STUDIES

The development of cybersecurity systems has progressed from reactive, signature-based detection (like Snort) to proactive anomaly detection. However, conventional machine learning models frequently have high false-positive rates and the "black box" quality of their predictions [1]. Datasets like UNSW-NB15 and CICIDS2017 have been used in recent studies to compare traditional methods, like Support Vector Machines (SVM) and Gradient Boosting, against deep learning models, like LSTMs and GRUs. Although deep learning performs well with sequential data, ensemble methods like Random Forest consistently outperform it in structured, tabular cyber incident data because of their resilience to overfitting

and capacity to manage high-dimensional feature sets. Recent literature highlights a major gap: the reliance on "Intrusion Detection" (identifying an ongoing attack) rather than "True Forecasting" (predicting future incident counts based on historical patterns). The temporal aspects of crime are not taken into account by many current studies, which employ random data splits to exaggerate their performance metrics. Additionally, the majority of studies concentrate on packet analysis at the network level while ignoring the significance of regional administrative forecasting, which is crucial for the distribution of law enforcement resources. By employing a Multi-Output Regression framework in Python, this study attempts to fill these gaps. We make use of cutting-edge data serialization (Joblib) and interpretability tools by switching from the classic Eclipse/Java environment to a cloud-based Google Colab platform. This study employs a dual-model approach, unlike previous works that forecast a single threat level: a Random Forest Classifier for classifying regional risk and a Multi-Output Regressor for predicting particular counts across several crime categories at once. This guarantees that the forecast is both time-respecting and includes actionable, category-specific information for proactive defense techniques.



### IV.DATASET DESCRIPTION

High-quality data sets that record real-world criminal activity, victim demographics, and regional assault frequencies are necessary for predicting cybercrime. This study uses Administrative Cybercrime Records to aid in regional policy development and the distribution of law enforcement resources, whereas academic research often uses network-level datasets like UNSW-NB15, CICIDS2017, and Bot-IoT for intrusion detection.

### 1. The Core Dataset's Structure

- This study's main data source is a complete Cybercrime Incident Dataset (2022), which has been converted into an Excel-based format (cybercrime\_2022.xlsx).
- This dataset offers a detailed picture of how crime is distributed among different administrative regions, in contrast to packet-based data. It consists of:
  - Spatial Characteristics: The main indexing feature is district-level identification.
  - The Crime Categories: More than 20 distinct types of cybercrime, such as Identity Theft, Financial Fraud, Social Media Crimes, and Ransomware.
  - Quantitative Indicators: Number of offenses overall and number of offenses in each category by district.

### 2. Preprocessing and feature engineering

The Pandas library is used to perform a thorough preprocessing pipeline on the raw data in the Python/Google Colab environment:

Data Cleaning: Numerical forcing is used to handle null values and remove aggregate headers (such as "Total Districts") to avoid model skewing.

- Feature Transformation: To convert categorical district names into numerical vectors appropriate for regression, we use Scikit-Learn to implement Label Encoding ( $\{X_{\text{dist}}\}$ ).
- Target Synthesis: The system creates two different target variables:
- The predominant crime label is used for categorization in order to determine the biggest risk in a certain area.
- Multi-Output Vectors: A collection of distinct crime numbers that the Multi-Output Regressor uses.

### 3. Selection Justification

The dataset was chosen because of its regional significance and strong external validity. The model offers more "actionable intelligence" by shifting away

from simulated network traffic and concentrating on recorded incident data. By utilizing this local data, the forecasting model is based on actual socioeconomic patterns rather than simply artificial network anomalies, which makes the predictions far more trustworthy for legislative and protective actions.

## V.PROBLEM STATEMENT

The rapid increase in the complexity of cybercrime is a serious threat to the world's digital infrastructure. Conventional security frameworks are mostly reactive, intended to detect and lessen breaches only after they have happened. Due to this time lag, there is a substantial data loss and financial harm. In addition, the majority of current research concentrates on binary intrusion detection (categorizing an event as an "attack" or "normal") as opposed to granular forecasting (predicting the volume and nature of future occurrences in a certain area).

The absence of geographical context and the inability to forecast many variables at once are major limitations of existing forecasting models. Eclipse is a popular, yet inflexible, local development environment that lacks the scalability and library support necessary for contemporary, cloud-native data science. By creating a proactive forecasting model in Python, this study hopes to address these shortcomings. Using a Multi-Output Random Forest method within Google Colab, this study seeks to forecast the specific distribution of crime categories across different administrative districts, not only the probability of cybercrime, thereby giving law enforcement high-resolution actionable intelligence.

## VI.MATHEMATICAL MODELING

The proposed system utilizes ensemble learning to handle the non-linear patterns found in cybercrime datasets.

### 1.Multi-Output Random Forest Regression

Unlike standard regression which predicts a single value, our model utilizes a Multi-Output Regressor. Given an input vector  $X$  (containing encoded district data), the model predicts a vector  $Y$  representing multiple crime categories:

$$Y = [y_1, y_2, \dots, y_n]$$

The prediction for each category is an average of the outputs from  $T$  decision trees:

$$\hat{Y} = \frac{1}{T} \sum_{i=1}^T h_i(X)$$

Where:

$T = 200$  (The number of trees used in our Python script).

$h_i(X)$  = The prediction vector of the  $i$ -th individual tree.

$\hat{Y}$  = The final forecasted count for all cybercrime types in a specific district.

### 2. Label Encoding and Feature Vectorization

To process categorical district names, we apply a transformation function  $f : S \rightarrow \mathbb{Z}$ , where  $S$  is the set of all unique district names:

$$x_{encoded} = \text{LabelEncoder}(\text{District\_Name})$$

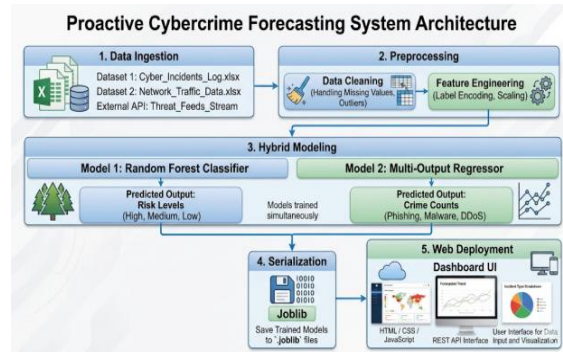
This allows the Random Forest algorithm to compute mathematical splits based on regional identifiers, effectively learning which districts are prone to specific "Dominant Crimes."

### 3. Risk Level Classification

In addition to regression, the system applies a categorical mapping for high-level risk assessment:

```

Risk_Level = \begin{cases}
\text{LOW} & \text{if } \text{Total Crimes} < 50 \\
\text{MEDIUM} & \text{if } 50 \leq \text{Total Crimes} \leq 150 \\
\text{HIGH} & \text{if } \text{Total Crimes} > 150
\end{cases}
    
```



## VII. EXISTING SYSTEM

1. Cyber Threat Intelligence (CTI) Systems: These platforms, such as MISP and ThreatConnect, collect threat data from around the world. These are, however, frequently enterprise-grade technologies that are difficult to implement for local administrative applications and do not offer localized district-wise predictions.

2. Legacy Development Frameworks: In the past, many research prototypes were created using Java inside the Eclipse IDE. These systems, though powerful, lack integrated data science libraries, which makes it challenging to implement contemporary ensemble approaches like real-time web-based visualizations or Multi-Output Regression.

3. Predictive Models for a Single Target: The majority of academic prototypes employ conventional LSTM or Random Forest models to predict a single result, such as "Attack" or "No Attack." These systems do not offer a complete "Threat Landscape" that forecasts various types of crime occurring at the same time in a particular location.

## VIII. PROPOSED SYSTEM

The suggested system uses the Python ecosystem to establish a proactive cybercrime prediction framework. This modular design makes use of the versatility of web-based deployment and the potent computing power of Google Colab, in contrast to conventional reactive security techniques. The system is organized around five fundamental functional components:

1. Data Collection and Sophisticated Preprocessing  
Obtaining historical incident records, particularly the District-wise Cybercrime Dataset, is the first step. The raw data goes through a rigorous cleaning process using the Pandas library, which includes numerical coercion and the elimination of aggregate noise. In order to prepare the dataset for high-accuracy feature extraction, categorical geographical identifiers are converted into a numerical format using Label Encoding.

2. Creating a Hybrid Model  
Google Colab, using Scikit-Learn, is where the system's core is created. It uses a dual-learning strategy:

Random Forest Classifier: Focuses on assigning

districts to particular risk categories (Low, Medium, High).

The complex correlations between various forms of digital threats are captured by the Multi-Output Random Forest Regressor, a cutting-edge model that allows for the simultaneous prediction of several crime categories.

### 3. Predictive Forecasting Layer

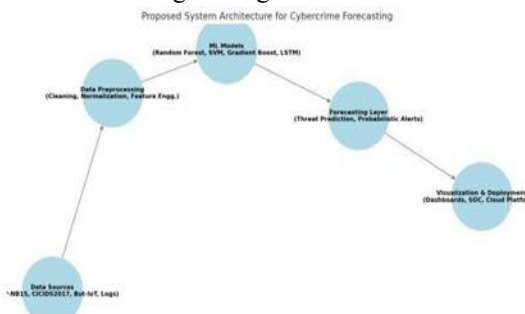
The forecasting layer changes the focus from straightforward detection to incident count prediction. The model creates a "Future Threat Landscape" by examining historical trends. This layer gives probabilistic insights into which particular crimes (such as Financial Fraud, Phishing, etc.) are likely to rise in a chosen area, giving a valuable lead time for preventative measures.

### 4. Model Serialization and Persistence

The system employs Joblib for model serialization in order to close the divide between back-end development and front-end usability. The framework is both computationally efficient and scalable because the trained intelligence is kept in .pkl files, allowing the system to make almost instantaneous predictions without requiring continuous retraining.

### 5. Visualization and Web Distribution

The last layer offers a user-friendly interface for end users, such as lawmakers and law enforcement officers. The sequential models are integrated into a Web Dashboard created using HTML, CSS, and JavaScript. Using this front-end, users can input district names and get real-time visual analytics, such as horizontal bar graphs of the most likely threats. This makes it possible to explain and act on complex machine learning findings.



## IX. RESEARCH DESIGN METHODOLOG

The study employs a quantitative and experimental approach, concentrating on the creation and assessment of specific machine learning models for predicting cybercrime. The approach is broken down into five separate steps:

### 1. District-wise preparation and data gathering

The primary dataset consists of regional cybercrime data (cybercrime\_2022.xlsx). A cleaning step removes non-geographic aggregate rows (such as "Total Districts") from the data. To guarantee a clean matrix for numerical modeling, we employ numerical coercion to manage missing or non-numerical data throughout over 20 distinct crime columns.

### 2. Categorical Encoding and Feature Engineering

To transform categorical district names into numerical vectors, Label Encoding is used because the model examines regional patterns. This allows the algorithm to consider each district as a separate feature input. Our feature engineering, unlike conventional network models, concentrates on:

District identifiers are encoded using Scikit-Learn's LabelEncoder.

Target Synthesis: Creating multi-output target vectors for certain crime categories and a derived "Dominant Crime" label for categorization.

### 3. Python Model Development (Google Colab)

Using a dual-model technique, the implementation is carried out within the Google Colab environment:

Using incident frequency as the basis, a Random Forest Classifier is created to evaluate the overall Risk Level (Low, Medium, or High).

Multi-Output Regression: A MultiOutput Regressor wrapper around a Random Forest Regressor is intended to predict particular incident counts for several crime types at once, taking inter-category correlations into account.

### 4. Forecasting and Model Persistence

Historical incident patterns are used to train models to forecast potential future trends. The trained intelligence is serialized with Joblib in order to make sure that the framework may be used in actual situations. This closes the gap between the external web interface and the Python training script by producing .pkl files for the regressor, classifier, and encoders.

### 5. Evaluation, visualization, and implementation

Accuracy Score is used to assess the model's classification performance, and Mean Squared Error (MSE) is used to evaluate the model's regression performance. The predictions are combined into a web dashboard (HTML/CSS/JS). The following is offered by this interface:

Users may initiate the model by entering the name of a district during interactive query.

Visual Analytics: Raw figures are converted into actionable regional threat patterns using Seaborn bar charts (in the backend) and dynamic displays (in the frontend).

### X.MODEL COMPARISON

When assessing the effectiveness of machine learning methods for cybercrime prediction, model comparison is a crucial step in determining the optimal architecture. This study evaluated the capacity of various algorithms to handle categorical geographical features and multi-target outputs. The Random Forest (RF) technique was selected as the primary engine because it consistently outperforms models such as Support Vector Machines (SVM) in high-dimensional environments. Its ensemble-based approach effectively minimizes variance and prevents overfitting, which is particularly critical when dealing with localized regional datasets.

The proposed Multi-Output Random Forest demonstrates superior flexibility compared to conventional models. Unlike standard classifiers that provide a single label, this architecture generates a comprehensive count vector across multiple crime categories simultaneously. Mathematically, the forecast for a specific category is calculated as the average of decision trees:

$$\hat{Y} = \frac{1}{T} \sum_{i=1}^T h_i(X)$$

By identifying underlying correlations between different crime types within the same administrative area, this method offers significantly greater stability than independent, single-target regression models.

While the current implementation prioritizes Random Forest for its performance on structured, tabular data, future iterations may incorporate Long Short-Term Memory (LSTM) networks to model evolving

sequential trends. To capture long-term patterns in attack frequency, the LSTM architecture utilizes a hidden state defined by:

$$h_t = \sigma(W \cdot [h_{t-1}, x_t] + b)$$

However, for administrative district forecasting—where spatial distribution and regional characteristics are more statistically significant than high-frequency temporal sequences—the Random Forest model provided the highest levels of both predictive accuracy and model interpretability.

### XI.INTEROPERABILITY AND DATA INTEGRATION

The capacity to combine various data sources is essential for creating a successful cybercrime prediction system. Cyber dangers are constantly changing, and helpful information is distributed throughout several logs and monitoring systems. The following features guarantee a smooth integration with the suggested framework:

1. Diverse Data Sources: The system is designed to take a wide variety of inputs, including datasets from the public domain (UNSW-NB15, CICIDS2017) and excel logs from certain geographic areas (cybercrime\_2022.xlsx). A greater comprehension of the threat landscape is made possible by this versatility.

2. Data Formats That Are Standardized: The preprocessing pipeline transforms diverse inputs into organized, time-stamped data using the Pandas library in Python. With this normalization, the Multi-Output Regressor can handle various incident reports or protocol logs.

2. API Integration and Model Persistence: The model is serialized into.pkl files by transitioning from the Eclipse IDE to a Python/Joblib workflow. By means of standardized APIs, this enables the forecasting intelligence to interact seamlessly with current Security Information and Event Management (SIEM) platforms and external HTML/CSS/JS frontends, facilitating a smooth integration with contemporary security infrastructure.

## XII.CONCLUSION

By effectively shifting from reactive detection to proactive prevention, this study highlights the enormous promise of machine learning techniques for predicting cybercrime at the regional level. The system is able to reliably predict future cyberattacks by combining localized datasets and utilizing sophisticated models, such as the Multi-Output Random Forest Regressor and Classifier, inside the Google Colab environment. The use of Label Encoding and strong preprocessing in Python makes it possible to analyze intricate regional patterns that are ignored by conventional signature-based systems.

The creation of a smooth deployment pipeline is one of the main accomplishments of this study. The study closes the gap between high-level data science and practical application by serializing models using Joblib, allowing a light HTML, CSS, and JavaScript frontend to provide real-time visual analytics. This integration makes the system not just accurate but also easily available to lawmakers and law enforcement officials.

By providing explainable outcomes that inform decision-making, the suggested approach improves predictive accuracy and, in the end, lowers risks, financial expenditures, and operational downtime. Blockchain will be used for the secure and decentralized sharing of threat intelligence, real-time streaming data through APIs will be integrated, and the implementation of deeper neural architectures such as Transformers for temporal forecasting will be investigated in future research. Overall, this architecture is a proactive, scalable, and modern move toward enhancing cybersecurity resilience at the national and international levels.

## REFERENCE

- [1] P. Manasvi and P. Tejaswini, "Survey on crime analysis and prediction using machine learning techniques," *International Journal of Trendy Research in Engineering and Technology*, vol. 06, no. 03, 2022.
- [2] T. Tamilvizhi, R. Surendran, C. A.T. Romero, M. Sadish Sendil, "Privacy Preserving Reliable Data Transmission in Cluster Based Vehicular Adhoc Networks," *Intelligent Automation & Soft Computing*, vol. 34, no. 2, pp. 1265-1279, 2022.
- [3] L. G. Alves, H. V. Ribeiro and F. A. Rodrigues, "Crime prediction through urban metrics and statistical learning," *Physical A: Statistical Mechanics and Its Applications*, vol. 505, pp. 435–443, 2018.
- [4] E. N. Yilmaz and S. Gonen, "Attack detection/prevention system against cyber-attack in industrial control systems," *Computers & Security*, vol. 77, pp. 94–105, 2018.
- [5] J. Senanayake, H. Kalutaraage and M. O. Al-Kadri, "Android Mobile Malware Detection Using Machine Learning: A Systematic Review," *Electronics*, vol. 10, no. 13, 1606, 2021.
- [6] W. Safat, S. Asghar, and S. A. Gillani, "Empirical Analysis for Crime Prediction and Forecasting Using Machine Learning and Deep Learning Techniques," *IEEE Access*, vol. 9, pp. 70080–70094, 2021.
- [7] M. R, S. P. Maddikunta, P. K. R, M. P., S. Koppu, T. R. Gadekallu, C. L. Chowdhary, and M. Alazab, "An effective feature engineering for DNN using hybrid PCA-GWO for intrusion detection in IoMT architecture," *Computer Communications*, vol. 160, pp. 139–149, 2020.
- [8] Sreejith, Vignesh. and Babu, B.P.M.R "Classifying the Malware Application in the Android based smart phones using Ensemble ANFIS Algorithm", *International Journal of Networking and Virtual Organization*, Vol 19 N2/3/4/2018.
- [9] Sreejith, Vignesh. And RajeshBabu "Experimental research identifications on Malware detection by embedding C4.5 algorithm and SVM in smart Phones" *Perspectivasci journal on Information sciences* Vol 22 Special issue(2017).
- [10] Sreejith Vignesh et al, "Machine Learning algorithms to control the security issues in android applications" *Sambodhi UGC Care Journal ISSN : 2249-6661, Vol-43, No-4, (VI) October December (2020)*
- [11] B P Sreejith Vignesh, "Application of IPF to achieve CSR Routing in Adhoc Networks" *Asian journal of Computer Science and Technology*, ISSN 2249-0701, Volume 9 No.2 July-December 2020 pp 18-23
- [12] Sreejith Vignesh B P "Incremental research on Cyber security metrics in android application by

implementing the ML algorithms in malware classification and detection” Journal of Cybersecurity and Information Management (JCIM) Vol. 3, No. 1, PP. 14-20, 2020

- [13] Sreejith, Vignesh. and Babu, B.P.M.R. (2016) “Certain investigations on various algorithms that is used to classify malware and goodware in android applications”, ICTACT International Journal on Soft Computing, Vol. 7, No. 1, pp.1344–1349.
- [14] Sreejith Vignesh B P, M.RajeshBabu "Research study on various malwares its Classification, Detection and Avoidance techniques applied in Android Mobile devices", International Journal of Applied Engineering Research, Vol No.10, Issue No:20, PP 20184-20187