# ONLINE PAYMENTS FRAUD DETECTION USING WITH MACHINE LEARNING:

To build an application that can detect the legitimacy of the transaction in real-time and increase the security to prevent fraud.

## By

*(Marri yashmitha)*

*(Manthina raja rishika)*

*(Kutagulla safa)*

*Guided by*

## *Prof.  Ms swetha raj*

A Dissertation Submitted to
SRI  VENKATESWARA  COLLEGE  OF
ENGINEERING AND TECHNOLOGY, An
Autonomous Institution affiliated to
'JNTU Ananthapur' in Partial Fulfilment of
the  Bachelor  of  Technology  branch  of
*Computer science and Engineering*

*May 2024*

# SRI VENKATESWARA COLLEGE OF ENGINEERING AND TECHNOLOGY

### R.V.S. Nagar Tirupathi Road, Andhra Pradesh– 517127

## Data Collection and Preprocessing for Fraud Detection:

Collection and Preprocessing Data collection and preprocessing
Data Collection and Preprocessing for Fraud Detection

 In the realm of fraud detection, accurate data collection and preprocessing are crucial steps in building effective models. The quality and relevance of the data directly impact the performance of *fraud detection* algorithms, making it vital to carefully consider various aspects of *data collection* and preprocessing. From selecting the right data sources to handling missing values and outliers, these steps play a pivotal role in uncovering fraudulent activities. This section delves into the intricacies of *data collection* and preprocessing, shedding light on various perspectives and providing insights on the best practices.

1. Identifying relevant data sources: The first step in data collection is to identify the most relevant data sources that can provide valuable insights into fraudulent activities. These sources can include transaction logs, customer profiles, device information, and external databases. By leveraging multiple data sources, organizations can gain a comprehensive understanding of fraudulent patterns and enhance the accuracy of their *fraud detection models*.

2. Handling missing values: Missing data is a common challenge in *fraud detection*, and it can significantly impact the performance of models if not

handled properly. There are several approaches to address missing values, such as *imputation techniques* like mean, median, or mode imputation, or more advanced methods like regression imputation or multiple imputations. The choice of imputation method depends on the nature of the data and the potential impact on the final model. For example, if a significant portion of a specific feature is missing, it might be better to exclude that feature altogether rather than imputing the *missing values*.

3. Dealing with outliers: Outliers are extreme values that deviate significantly from the normal distribution of data. These outliers can arise due to errors, anomalies, or fraudulent activities. Handling outliers is crucial to ensure the robustness of *fraud detection models*. One approach is to remove outliers based on *statistical techniques* such as z-score or interquartile range. Another option is to use *robust algorithms*, like the Isolation Forest or *Local Outlier Factor*, that can effectively identify and handle outliers without compromising the overall performance of the model.

4. Feature engineering and selection: Feature engineering involves transforming raw data into meaningful features that capture the underlying patterns of fraudulent activities. Techniques such as binning, one-hot encoding, scaling, and creating interaction terms can enhance the predictive power of the features. Additionally, feature selection methods like correlation analysis, information gain, or regularization techniques can help identify the most relevant features for fraud detection models. Striking the right balance

between feature engineering and feature selection is essential to avoid overfitting and improve the model's generalizability.

5. Balancing the data: Imbalanced datasets, where the number of fraudulent instances is significantly lower than the legitimate ones, pose a challenge for fraud detection models. Ignoring the imbalance can lead to biased models that fail to detect fraudulent activities effectively. To address this, various techniques can be employed, including oversampling the minority class (fraudulent instances) through methods like SMOTE or undersampling the majority class (legitimate instances). Alternatively, ensemble techniques like boosting or bagging can be used to combine multiple models trained on balanced subsets of the data. The choice of balancing technique depends on the specific dataset and the trade-off between recall (*fraud detection* rate) and precision (accuracy of *fraud detection*).

Overall, data collection and preprocessing are intricate processes that require careful consideration in fraud detection. By selecting relevant data sources, handling missing values and outliers, performing effective feature engineering and selection, and balancing the data, organizations can build robust fraud detection models that accurately identify fraudulent activities while minimizing false positives. These steps, when executed with a comprehensive understanding of the data and the specific requirements of the problem, lay a solid foundation for successful *fraud detection* systems.