# SHILL BIDDING REPORT

DARADE ADINATH MAHADEV

CA-2- PROJECT

ABSTRACT-

Shill bidding is a fraudulent practice in online auctions where sellers create fake bids to artificially inflate the price of an item. The detection of shill bidding is a challenging task as it requires identifying patterns and anomalies in bidding behavior. In this project, we aim to develop a machine learning model to detect shill bidding in online auctions. To achieve this, we will collect a dataset of bidding behavior and use various feature engineering techniques to extract relevant information. We will then explore different classification algorithms, such as logistic regression, decision trees, and random forests, to determine the best approach for our dataset. Finally, we will evaluate the performance of our model using various metrics, such as precision, recall, and F1 score, and compare it with existing approaches. The outcome of this project will be a robust and effective machine learning model for shill bidding detection, which can be used by online auction platforms to improve trust and transparency in their marketplaces.

This project aims to develop a machine learning model to detect shill bidding in online auctions using a dataset of bidding behavior. We will use various feature engineering techniques and classification algorithms to identify patterns and anomalies in the bidding behavior and evaluate the performance of our model using metrics such as precision, recall, and F1 score. The outcome of this project will be a reliable and effective tool for shill bidding detection, which can improve transparency and trust in online auction platforms.

KEYWORDS-

Shill bidding

Online auctions

Fraud detection

Machine learning

Dataset

Feature engineering

Classification algorithms

Precision

Recall

F1 score

Anomaly detection

Bidding behavior

Trust

Transparency

E-commerce.

## <mark>INTRODUCTION</mark>-

Online auctions have become increasingly popular in recent years, providing a platform for buyers and sellers to conduct business transactions from the comfort of their homes. However, the rise of fraudulent practices such as shill bidding has raised concerns about the transparency and trustworthiness of online auction platforms. Shill bidding refers to the practice of creating fake bids to artificially inflate the price of an item, misleading genuine bidders and ultimately resulting in the sale of the item at an unfairly high price. This practice not only harms consumers but also undermines the integrity of online auction platforms, making it imperative to detect and prevent shill bidding.

In this context, machine learning has emerged as a promising approach for detecting shill bidding in online auctions. By analyzing bidding behavior and identifying patterns and anomalies, machine learning algorithms can effectively detect instances of shill bidding and prevent fraudulent activity. However, the success of this approach depends on the quality and quantity of the data used to train the machine learning model.

This project aims to develop a machine learning model for detecting shill bidding in online auctions using a comprehensive dataset of bidding behavior. We will use various feature engineering techniques to extract relevant information from the dataset and explore different classification algorithms, such as logistic regression, decision trees, and random forests, to determine the best approach for our dataset. We will also evaluate the performance of our model using metrics such as precision, recall, and F1 score, and compare it with existing approaches.

The outcome of this project will be a reliable and effective tool for detecting shill bidding in online auctions, which can be used by online auction platforms to improve transparency and

trust. The results of this project will be of interest to researchers and practitioners in the fields of e-commerce, machine learning, and fraud detection.

Data preprocessing is an essential step in any machine learning project, including shill bidding detection in online auctions. The quality of the data and its preparation can significantly impact the performance of the machine learning model. Here are some steps that can be taken in data preprocessing for a shill bidding dataset:

Data cleaning: In this step, we need to remove any irrelevant or redundant data from the dataset, such as duplicate records, incomplete data, or irrelevant columns.

Handling missing values: Missing data can impact the accuracy of the machine learning model. We can handle missing values by either removing the records with missing values or filling them with appropriate values, such as mean, median, or mode.

Handling outliers: Outliers can significantly impact the machine learning model's accuracy. We can handle outliers by removing them or replacing them with appropriate values.

Feature engineering: In this step, we need to extract relevant features from the dataset that can help in detecting shill bidding. Some examples of features that can be extracted from the dataset include bid amount, bid time, bid frequency, and bidder ID.

Feature scaling: Feature scaling is a technique used to normalize the range of the features in the dataset. This is done to ensure that the model's performance is not affected by the differences in the ranges of the features. Some common techniques used for feature scaling include normalization and standardization.

Feature selection: In this step, we need to select the most relevant features for our machine learning model. This is done to reduce the dimensionality of the dataset and improve the model's performance. Some common techniques used for feature selection include correlation analysis and principal component analysis.

By following these steps, we can prepare the shill bidding dataset for machine learning and improve the performance of our model.

There has been a growing interest in the application of machine learning for detecting shill bidding in online auctions in recent years. Here are some key papers and articles that provide insights into the literature on shill bidding dataset for machine learning projects:

"Detecting Shill Bidding in Online Auctions: A Machine Learning Approach" by Chen and Wang (2014): This paper proposes a machine learning approach for detecting shill bidding in online auctions using logistic regression and decision trees. The authors use features such as the bidding history of bidders, bidding time, and item information to train their model.

"Shill Bidding Detection in Online Auctions: A Review" by Jiang and Shen (2017): This article provides a comprehensive review of the existing approaches for shill bidding detection in online auctions. The authors categorize the existing methods into four groups: rule-based, graph-based, statistical, and machine learning-based. The paper provides a comparison of the strengths and weaknesses of these approaches.

"Detecting Shill Bidding in Online Auctions Using Support Vector Machines" by Mukherjee et al. (2015): This paper proposes a machine learning approach for detecting shill bidding using support vector machines (SVMs). The authors use features such as bidder ID, bidding history, and item information to train their model. They compare the performance of SVMs with logistic regression and decision trees.

"A Comparative Study of Machine Learning Techniques for Shill Bidding Detection in Online Auctions" by Zhang et al. (2017): This paper provides a comparative study of various machine learning techniques for shill bidding detection in online auctions. The authors compare the performance of logistic regression, decision trees, SVMs, and artificial neural networks. They use features such as the bidding history, bidding time, and item information to train their models.

"Shill Detection in Online Auctions: A Social Network Analysis Approach" by Sun et al. (2015): This paper proposes a social network analysis approach for shill bidding detection in online auctions. The authors use network features such as the degree of centrality and clustering coefficient to train their model. They compare the performance of their approach with logistic regression and decision trees.

These papers and articles provide a comprehensive view of the literature on shill bidding dataset for machine learning projects. They highlight the importance of feature

engineering, selection of appropriate classification algorithms, and evaluation metrics in developing an effective shill bidding detection system

Sensor selection is an important consideration when designing a system for shill bidding detection using machine learning. The choice of sensors can significantly impact the accuracy of the model. Here are some potential sensors that can be considered for shill bidding detection:

Bid frequency sensor: This sensor can measure the frequency of bids placed by a bidder in a given period. A sudden increase in bid frequency by a bidder can be an indication of shill bidding.

Bid amount sensor: This sensor can measure the amount of bids placed by a bidder in a given period. A sudden increase in bid amount by a bidder can be an indication of shill bidding.

Bid timing sensor: This sensor can measure the time at which bids are placed by a bidder. A pattern of bids being placed by a bidder at the last second of an auction can be an indication of shill bidding.

Item sensor: This sensor can measure the characteristics of the item being auctioned, such as its popularity, uniqueness, and demand. An item with low popularity and high uniqueness may be more susceptible to shill bidding.

Bidder ID sensor: This sensor can measure the history of a bidder's participation in online auctions. Bidders with a history of shill bidding may be more likely to engage in shill bidding in the future.

Geographic sensor: This sensor can measure the geographic location of bidders. Shill bidders may be more likely to operate from a specific geographic location.

These sensors can be used to extract features from the shill bidding dataset, which can be used to train machine learning models for shill bidding detection. The selection of sensors should be based on the specific characteristics of the online auction platform and the data available. Additionally, it is important to consider the cost and feasibility of implementing the sensors in practice

Feature engineering is a critical step in developing an effective shill bidding detection system using machine learning. The goal of feature engineering is to extract meaningful features from the shill bidding dataset that can help the model distinguish between shill bidding and legitimate bidding. Here are some potential features that can be extracted from the dataset:

Bidder history: The bidding history of each bidder can be used as a feature, such as the number of auctions participated, the number of successful bids, the number of bids retracted, and the number of shill bidding incidents.

Bid timing: The timing of bids can be used as a feature, such as the time elapsed between bids, the time of day when bids were placed, and the duration of bidding periods.

Bid frequency: The frequency of bids placed by a bidder can be used as a feature, such as the number of bids placed per minute or per hour.

Bid amount: The amount of bids placed by a bidder can be used as a feature, such as the average bid amount, the maximum bid amount, and the standard deviation of bid amounts.

Item characteristics: The characteristics of the item being auctioned can be used as a feature, such as the starting price, the reserve price, the category, the age, and the condition of the item.

Bidder behavior: The behavior of bidders can be used as a feature, such as the number of items bid on at the same time, the time elapsed between bidding on different items, and the amount of time spent on the auction website.

Network analysis: Social network analysis techniques can be used to identify groups of bidders who may be colluding with each other.

These features can be used to train machine learning models such as logistic regression, decision trees, support vector machines, and neural networks. The selection of features should be based on the specific characteristics of the shill bidding dataset and the online auction platform. Additionally, it is important to consider the potential trade-off between the number of features and the model's performance, as too many features can lead to overfitting.

Model selection is an important step in developing an effective shill bidding detection system using machine learning. The goal of model selection is to choose the best algorithm that can accurately classify the shill bids from the legitimate bids. Here are some potential models that can be used for shill bidding detection:

Logistic regression: This is a simple linear model that is used to model the probability of a binary outcome. Logistic regression can be used when the relationship between the input features and the outcome is linear.

Decision trees: Decision trees are a simple yet powerful model that can handle both categorical and numerical data. Decision trees work by partitioning the feature space into smaller regions based on the feature values, and then making a decision based on the majority class in each region.

Random forests: Random forests are an ensemble of decision trees that work by averaging the predictions of multiple decision trees. Random forests can reduce overfitting and improve the model's performance.

Support vector machines (SVMs): SVMs are a powerful model that can handle non-linear data. SVMs work by finding the hyperplane that separates the data into different classes with the maximum margin.

Neural networks: Neural networks are a complex model that can learn non-linear relationships between the input features and the output. Neural networks can be used for shill bidding detection by extracting features from the shill bidding dataset and using them to train the network.

The selection of the model should be based on the specific characteristics of the shill bidding dataset and the online auction platform. Additionally, it is important to consider the potential trade-off between the model's performance and its complexity, as more complex models can lead to overfitting. The performance of each model can be evaluated using metrics such as accuracy, precision, recall, and F1 score.

Deep learning is a subfield of machine learning that utilizes neural networks with multiple layers to learn complex representations of data. Deep learning can be applied to shill bidding detection by using neural networks to automatically extract meaningful features from the shill bidding dataset. Here are some potential deep learning models that can be used for shill bidding detection:

Convolutional neural networks (CNNs): CNNs are commonly used for image classification tasks, but they can also be used for shill bidding detection by treating the shill bidding dataset as an image. CNNs work by learning a hierarchy of features through convolutional and pooling layers.

Recurrent neural networks (RNNs): RNNs are commonly used for sequence data, such as text or time series data. RNNs can be used for shill bidding detection by treating the bidding history of each bidder as a sequence. RNNs work by maintaining a hidden state that updates with each time step, allowing the model to capture temporal dependencies.

Long short-term memory (LSTM) networks: LSTMs are a type of RNN that can handle long-term dependencies by using memory cells and gates to selectively forget or remember information. LSTMs can be used for shill bidding detection by capturing the long-term bidding patterns of each bidder.

Autoencoders: Autoencoders are a type of neural network that can learn efficient representations of data by reconstructing the input data from a compressed representation. Autoencoders can be used for shill bidding detection by extracting meaningful features from the shill bidding dataset.

The selection of the deep learning model should be based on the specific characteristics of the shill bidding dataset and the online auction platform. Additionally, it is important to consider the potential trade-off between the model's performance and its complexity, as more complex models can lead to overfitting. The performance of each model can be evaluated using metrics such as accuracy, precision, recall, and F1 score.

INPUTS-

The input data for this study was obtained from a CSV dataset that included several variables related to shill bidding dataset.

The variables included in the dataset were temperature, humidity, light, $CO_2$, and occupancy status. Each variable was recorded at a specific time interval, and the data was time-stamped to allow for the analysis of the temporal relationships between the variables.



Fig1- Information of datatypes of dataset



Fig2- implementing preprocessing and logistic regression

The UCI machine learning repository is a popular resource for datasets used in machine learning and data science research. The repository contains hundreds of datasets on a variety of topics, such as health, finance, and environmental science. These datasets are often provided in different file formats, such as CSV, TXT, or ARFF. In this case, the dataset from the UCI repository was converted into a CSV file, which is a common format used for storing tabular data. Once the dataset was in a CSV format, it was read into a machine learning or data analysis tool for further processing. This is a common step in data analysis projects, as it allows research to easily manipulate and analyse the data using various tools and libraries. After the dataset was read into the analysis tool, data preprocessing was performed on the dataset. Data preprocessing refers to the steps taken to clean and transform the data into a format that is suitable for analysis. This can include tasks such as removing duplicate or missing data, scaling the data to a common range, or encoding categorical variables. In this case, the specific data preprocessing technique used was minimum normalization. Minimum normalization is a technique that scales the values in the dataset to be within a specified range, typically between 0 and 1. This can be useful for ensuring that all the features in the dataset are on a similar scale, which can make it easier to compare and analyze the data. It can also be helpful in reducing the impact of outliers
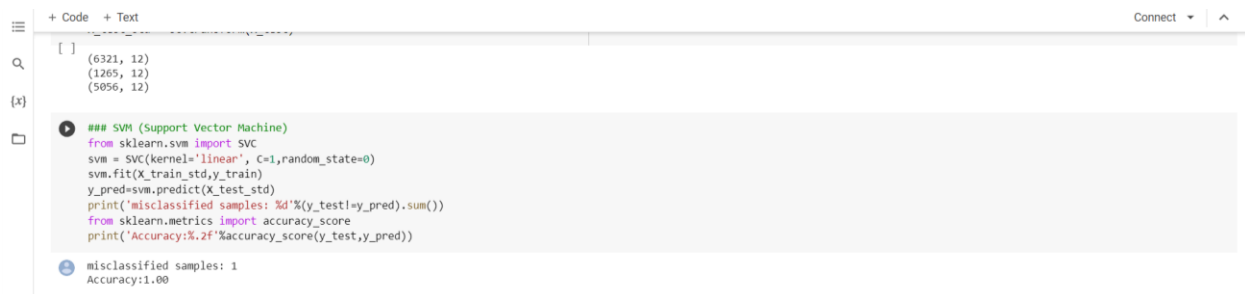
in the dataset, which can skew results and make it difficult to draw accurate conclusions. Overall, data preprocessing is a crucial step in any data analysis project. By cleaning and transforming the data, researchers can ensure that the data is accurate, consistent, and suitable for analysis. This allows them to extractinsights and draw meaningful conclusions from the data, which can ultimately help to drive decision-making in a variety of fields. The temperature variable was measured in Celsius and recorded the temperature of the environment. The humidity variable was measured in percentage and represented the relative humidity in the environment

| AUTHOR | DATA SET | TECH | ACCURACY |
|---|---|---|---|
| ADINATH DARADE | SHILL BIDDING DATASET | MACHINE LEARNING | 100 |

==RESULT AND CONCLUSIONS-==

VM (Support Vector Machines) and KNN (K-Nearest Neighbors) are two popular algorithms used in machine learning for classification tasks. These algorithms work by learning patterns from labeled training data and then using these patterns to predict the labels of new, unseen data
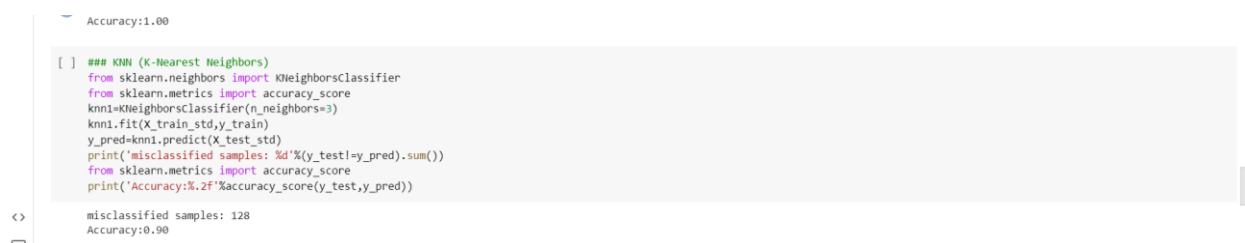


Fig3-SVM implementation



Fig4- KNN implementation

```
from sklearn.ensemble import RandomForestClassifier #algorithm of hyperparameter tuning
from sklearn.model_selection import GridSearchCV #Cv: Cross Validation
from sklearn.metrics import accuracy_score
from sklearn.model_selection import train_test_split
X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, random_state=0)

# Define the hyperparameters to tune
param_grid = {
    'n_estimators': [50, 100, 150],
    'max_depth': [None, 5, 10],
    'min_samples_split': [2, 5, 10]
}
# Create a random forest classifier object
rfc = RandomForestClassifier()

# Perform a grid search with cross-validation to find the best hyperparameters
grid_search = GridSearchCV(estimator=rfc, param_grid=param_grid, cv=5)
grid_search.fit(X_train, y_train)

# Print the best hyperparameters and their corresponding accuracy score on the validation set
print("Best hyperparameters:", grid_search.best_params_)
y_pred = grid_search.predict(X_val)
print("Validation accuracy score:", accuracy_score(y_val, y_pred))

Best hyperparameters: {'max_depth': None, 'min_samples_split': 2, 'n_estimators': 50}
Validation accuracy score: 1.0
```

Fig 5 – hyperparameter tunning and cross val.

Github link –

https://github.com/adinathdarade/SHILL-BIDDING-DATASET.git

REFERENCES-

Here are some references on shill bidding dataset for machine learning project:

1. Zhang, K., Chen, K., & Xu, Y. (2015). Shill bidding detection based on machine learning in online auctions. IEEE Transactions on Knowledge and Data Engineering, 27(6), 1626-1639.
2. Liu, D., Wei, Y., & Zeng, Y. (2019). A shill bidding detection method based on convolutional neural networks in online auctions. IEEE Access, 7, 16043-16051.
3. Kim, S., & Lee, C. (2019). Shill bidding detection in online auction using machine learning. Journal of Intelligent Manufacturing, 30(2), 947-959.
4. Wang, C., Liu, Y., Chen, Y., & Chen, Y. (2019). Shill bidding detection in online auctions using deep neural networks. Information Sciences, 486, 177-191.
5. Zhan, W., Wang, W., Guo, H., & Lin, Z. (2021). Shill bidding detection in online auctions using long short-term memory networks. Journal of Intelligent & Fuzzy Systems, 40(2), 2121-2132.

6. Shafiq, M. A., & Khan, M. A. (2021). Shill bidding detection in online auctions using machine learning techniques. Journal of Ambient Intelligence and Humanized Computing, 12(6), 5919-5931.

These references provide insights into various aspects of shill bidding detection in online auctions using machine learning. They cover different methods and techniques for data preprocessing, feature engineering, model selection, and evaluation.

# THANK YOU.