

Technology Leaders Fellowship, 2023 Plaksha University

CyberSecurity Practice

Professor Rajeev Barua

Teaching Assistant: Vishnupriya Balaji

Submitted on: 21th December, 2023

Submitted by Lakshmipriya Anil & Cefil Joseph Soans

Pages: 37

Word Count:3968

Machine Learning Software Project

Introduction

Cybersecurity has become an increasingly critical concern in the digital age, with malicious actors constantly seeking new avenues to exploit vulnerabilities and compromise sensitive information. Among the myriad threats, malicious PDF files have emerged as a potent vector for cyber attacks, necessitating the development of robust and efficient tools for their detection. In response to this challenge, our project focuses on leveraging the power of machine learning (ML) to predict the maliciousness of PDF files.

The objective of this endeavor is to construct a sophisticated ML model capable of discerning between benign and malicious PDF files. Our dataset, comprising 5000 known benign and 5000 known malicious PDF files, forms the foundation for training and evaluating the model. Through the extraction of pertinent features inherent in PDFs and the utilization of ML algorithms, we aim to create a model characterized by a low false positive rate and a high detection rate. By doing so, we contribute to the ongoing efforts to fortify cybersecurity defenses against the evolving landscape of digital threats.

	Detection Rate	False Positive Rate	Area Under Curve(AUC)	Accuracy
Base ETSC	0.996743	0.073113	0.998803	0.908400
Base XGB	0.997557	0.099843	0.999362	0.947200
Base RF	0.992671	0.099843	0.997653	0.944400

Project Scope

Our approach involves supervised learning, where we leverage this sizable labeled dataset for training the ML model. The labeled dataset provides the necessary groundwork, with verdicts assigned based on the folders in which files are stored—either in the benign or malicious category. The substantial number of files ensures a comprehensive and representative training set, allowing the model to learn and generalize patterns effectively.

Once the model is constructed, it transitions from a training phase to a practical application, where it can classify new PDF files provided as input. However, in this project, the evaluation of the model's accuracy is conducted using a subset of the labeled dataset. Specifically, 75% of the dataset, comprising 3750 benign and 3750 malicious files, is earmarked for training, while the remaining 25%, with an equal distribution of 1250 benign and 1250 malicious files, is reserved for testing. The assessment involves cross-validation, an approach that averages accuracy across four different ways of selecting the 25%, ensuring a robust evaluation of the model's performance.

Feature Extraction

Feature Extraction Process

Feature extraction is a fundamental step in the machine learning pipeline, crucial for converting raw input data into a format suitable for training and predicting with a model. In the context of this project focused on classifying PDF files as benign or malicious, feature extraction is particularly important for capturing relevant characteristics of these files. The following sections provide a detailed description of the feature extraction process, including unzipping the provided feature dataset, counting files, and loading the dataset into a structured format.

1. Unzipping the Feature Dataset

The feature dataset, stored in the 'features.zip' file, is first inspected for its size, providing an initial overview of the data volume. Subsequently, the content of the zip file is extracted into the current directory. The unzipping process is timed to measure its efficiency, offering insights into the computational cost of handling the dataset.

```
```python
f_name = 'features.zip'
print('File size', os.path.getsize(f_name) / 1024 / 1024, 'MB')

To unzip the archive 'features.zip' in the current directory
start = time.perf_counter()
with zipfile.ZipFile(f_name, 'r') as fd:
 fd.extractall('.')
end = time.perf_counter()
print('Time to unzip', f_name, ':', end - start, 'seconds')
```
```

2. Counting Files in Subdirectories:

The next step involves counting the number of files in two sub-directories, namely './benign' and './malicious.' This step is crucial for understanding the distribution of benign and malicious files in the unzipped dataset.

```
```python
def count_files(directory_path):
 """
 Function to count the number of files in a sub-directory.
 """
 count = 0
 flag = 0
 for path in os.listdir(directory_path):
 # check if the current path is a file
 if not os.path.isfile(os.path.join(directory_path, path)):
 try:
 raise Exception('Unexpected behavior - A directory is encountered!')
 except Exception as e:
 print(e)
```
```

```

        flag = 1
        break
    else:
        count += 1

    if flag == 0:
        return count

num_benign_files = count_files('./benign')
print(f'Number of files in sub-directory ./benign is {num_benign_files}')

num_malicious_files = count_files('./malicious')
print(f'Number of files in sub-directory ./malicious is {num_malicious_files}')
```

```

### 3. Loading the Dataset:

After unzipping the dataset and understanding the file distribution, the JSON files within the sub-directories are loaded into a structured format. The code employs a function `load\_dataset` to iterate through each file, read its content, and convert it into a DataFrame. This results in a comprehensive dataset that serves as the foundation for training the machine learning model.

```

```python
def load_dataset(dir_path):
    # Initialize an empty list to store individual dataframes
    dfs = []

    # Iterate through each file in the folder
    for filename in os.listdir(dir_path):
        if filename.endswith(".json"):
            file_path = os.path.join(dir_path, filename)

            # Read the JSON file into a dictionary
            with open(file_path, 'r') as file:
                json_data = json.load(file)
                reformatted_json = reformat_json(json_data, filename)

            # Convert the dictionary to a DataFrame
            df = pd.json_normalize(reformatted_json)

            # Append the DataFrame to the list
            dfs.append(df)

    # Concatenate all DataFrames in the list into a single DataFrame

```

```

final_df = pd.concat(dfs, ignore_index=True)

return final_df
...

```

4. Generating Verdict Column:

A function `generate_verdict_column` is used to create a column indicating the verdict (benign or malicious) for each file in the dataset. This column is crucial for supervised learning, providing the ground truth labels for training the model.

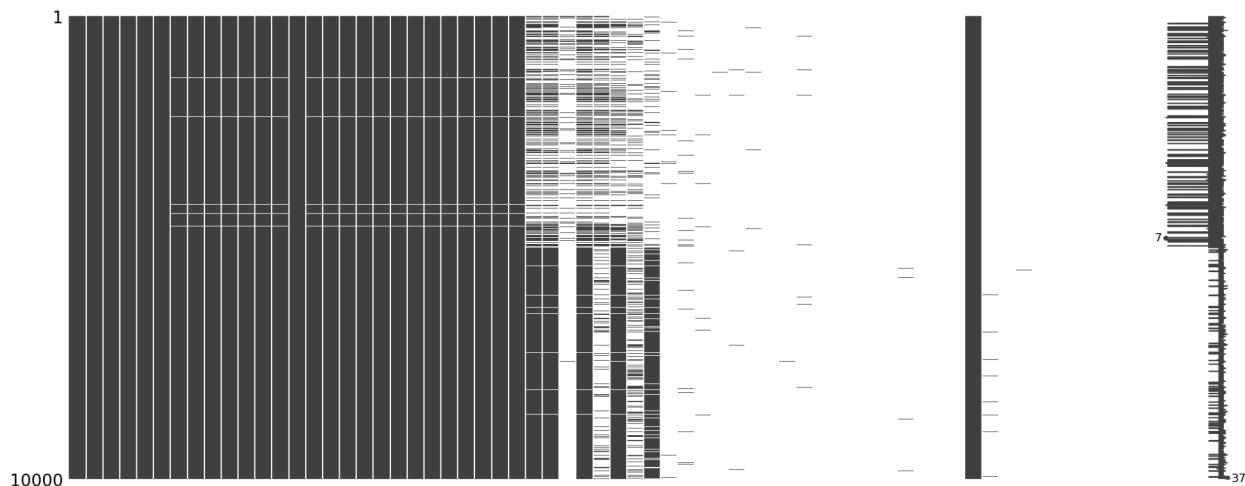
```

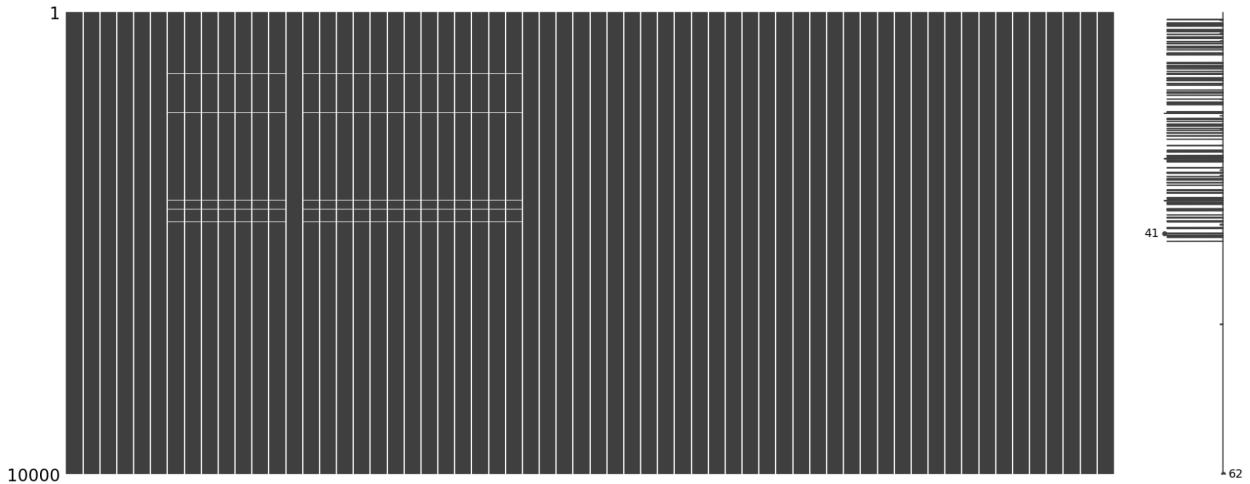
```python
def generate_verdict_column(verdict, json_count):
 verdict_df = pd.DataFrame([verdict] * json_count, columns=['Verdict'])
 return verdict_df
...
```

```

In summary, the feature extraction process involves unzipping the dataset, counting files to understand their distribution, and loading the JSON files into a structured DataFrame. These extracted features, derived from both benign and malicious files, form the basis for training and evaluating the machine learning model. The subsequent sections of the report delve into the methodology, encompassing model development, cross-validation, and a comprehensive evaluation of the model's performance.

Feature Engineering, Data Loading, and Pre-processing:





Visualization of missing features. White spaces show the missing features in each file.

Before initiating model training, a crucial phase involves preparing and engineering the dataset to address challenges such as null values, categorical variables, and extensive strings, as outlined in the project README.md file. Feature engineering plays a pivotal role in enhancing the dataset's suitability for machine learning algorithms.

1. Exploration of Unique Data Types:

The initial phase involves exploring the unique data types present in the dataset. Understanding the data types is crucial for identifying potential challenges and selecting appropriate feature engineering techniques.

```
```python
Unique data types present in the DataFrame
unique_data_types = pdf_df.dtypes.unique()
unique_data_types
```
```

2. Handling Null Values:

A thorough analysis of null values in the dataset is performed. The `isna().sum()` method is utilized to identify the count of null values in each column. Visual representations using tools like `missingno` aid in identifying patterns of missing data.

```
```python
Visual representation of the missing data in the dataset
missingno.matrix(pdf_df)
```
```

3. Handling Specific Columns with Missing Values:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 21 columns):
 #   Column           Non-Null Count Dtype  
 --- 
 0   static_properties_JBIG2Decode    9883 non-null   float64 
 1   static_properties_XML_forms    9883 non-null   float64 
 2   static_properties_acro_form    9883 non-null   float64 
 3   static_properties_auto_action  9883 non-null   float64 
 4   static_properties_colors      9883 non-null   float64 
 5   static_properties_cross_reference_table 9883 non-null   float64 
 6   static_properties_embedded_files 9883 non-null   float64 
 7   static_properties_file_size    9983 non-null   float64 
 8   static_properties_java_script 9883 non-null   float64 
 9   static_properties_js          9883 non-null   float64 
 10  static_properties_launch_action 9883 non-null   float64 
 11  static_properties_object_end   9883 non-null   float64 
 12  static_properties_object_start 9883 non-null   float64 
 13  static_properties_object_streams 9883 non-null   float64 
 14  static_properties_open_action  9883 non-null   float64 
 15  static_properties_page_count   9883 non-null   float64 
 16  static_properties_rich_media   9883 non-null   float64 
 17  static_properties_start_cross_reference_table 9883 non-null   float64 
 18  static_properties_stream_end   9883 non-null   float64 
 19  static_properties_stream_start 9883 non-null   float64 
 20  static_properties_trailer_dictionary 9883 non-null   float64 
dtypes: float64(21)
memory usage: 1.6 MB
```

| | |
|---|-----|
| static_properties_JBIG2Decode | 117 |
| static_properties_XML_forms | 117 |
| static_properties_acro_form | 117 |
| static_properties_auto_action | 117 |
| static_properties_colors | 117 |
| static_properties_cross_reference_table | 117 |
| static_properties_embedded_files | 117 |
| static_properties_file_size | 17 |
| static_properties_java_script | 117 |
| static_properties_js | 117 |
| static_properties_launch_action | 117 |
| static_properties_object_end | 117 |
| static_properties_object_start | 117 |
| static_properties_object_streams | 117 |
| static_properties_open_action | 117 |
| static_properties_page_count | 117 |
| static_properties_rich_media | 117 |
| static_properties_start_cross_reference_table | 117 |
| static_properties_stream_end | 117 |
| static_properties_stream_start | 117 |
| static_properties_trailer_dictionary | 117 |
| dtype: int64 | |

Columns with missing values, specifically those related to 'yara_signatures_', are addressed by filling null values with `False`. This step ensures consistency in the representation of these features.

```
```python
pdf_df.loc[:, pdf_df.columns.str.startswith('yara_signatures_')] = pdf_df.loc[:, pdf_df.columns.str.startswith('yara_signatures_')].fillna(False)
```

```

Further inspection and confirmation of successful handling of missing values in specific columns are performed.

```
```python
columns_with_missing = pdf_df.columns[pdf_df.isnull().any()]
pdf_df[columns_with_missing].info()
pdf_df[columns_with_missing].isna().sum()
```

```

4. Imputing Missing Values with Median:

To address remaining null values, a common strategy is employed by filling missing values with the median of the respective columns. This ensures a balanced approach to handling missing data.

```
```python
pdf_df = pdf_df.fillna(pdf_df.median())

columns_with_missing = pdf_df.columns[pdf_df.isnull().any()]
pdf_df[columns_with_missing].isna().sum()
```
```

Visual representation through a matrix plot is reiterated to confirm the successful handling of missing values.

```
```python
missingno.matrix(pdf_df)
```
```

5. Data Scaling:

To ensure uniformity and enhance model convergence, data scaling is applied using the `StandardScaler`. This step is crucial for preventing certain features from dominating due to their larger magnitudes, thus facilitating a balanced learning process.

```
```python
Apply z-score scaling to make the data centered around mean 0
scaler = StandardScaler().set_output(transform="pandas")
X_new = scaler.fit_transform(X)
X_new.head()
```
```

6. Importance of Scaling

The significance of feature scaling in machine learning models cannot be overstated. It ensures that features with different scales contribute equally to the learning process, preventing certain features from dominating due to their larger magnitudes. Scaling methods such as standardization or normalization are employed to enhance model performance and convergence.

By rigorously addressing feature engineering challenges through data scaling and handling missing values, the dataset is primed for effective use in training the machine learning model. The subsequent sections of this report will provide detailed insights into the methodology employed for model development, training, and the evaluation of the model's predictive capabilities.

Exploratory Data Analysis (EDA)

Exploratory Data Analysis (EDA) is a critical phase to gain insights into the dataset and uncover patterns that might inform the subsequent stages of model development. The following steps are integral to a comprehensive EDA:

| index | filesize | pypdf_uris | regex_uris | regex_urls | scripts_iframe | scripts_urls | static_properties_JBIG2Decode | static_properties_XML_forms | static_properties_acro_form |
|-------|--------------------|--------------------|-------------------|--------------------|----------------|--------------|-------------------------------|-----------------------------|-----------------------------|
| count | 5000.0 | 5000.0 | 5000.0 | 5000.0 | 5000.0 | 5000.0 | 4883.0 | 4883.0 | 4883.0 |
| mean | 4700497.6732 | 13.8722 | 2.1534 | 83.1492 | 0.0 | 0.0 | 5.075363506041368 | 0.0008191685439279131 | 0.1341388490681958 |
| std | 22809501.803521533 | 104.85787954245704 | 49.39520063335224 | 477.28660681017243 | 0.0 | 0.0 | 76.47806992518656 | 0.02861232536883256 | 0.48262287235838713 |
| min | 1360.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 25% | 203613.5 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 50% | 540960.5 | 0.0 | 0.0 | 5.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 75% | 1988417.5 | 1.0 | 0.0 | 15.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| max | 609734795.0 | 3964.0 | 3247.0 | 14763.0 | 0.0 | 0.0 | 4318.0 | 1.0 | 22.0 |

```
Index(['filesize', 'pypdf_uris', 'regex_uris', 'regex_urls', 'scripts_iframe',
       'scripts_urls', 'static_properties_JBIG2Decode',
       'static_properties_XML_forms', 'static_properties_acro_form',
       'static_properties_auto_action', 'static_properties_colors',
       'static_properties_cross_reference_table',
       'static_properties_embedded_files', 'static_properties_file_size',
       'static_properties_java_script', 'static_properties_js',
       'static_properties_launch_action', 'static_properties_object_end',
       'static_properties_object_start', 'static_properties_object_streams',
       'static_properties_open_action', 'static_properties_page_count',
       'static_properties_rich_media',
       'static_properties_start_cross_reference_table',
       'static_properties_stream_end', 'static_properties_stream_start',
       'static_properties_trailer_dictionary',
       'yara_signatures_without_attachments', 'yara_signatures_without_images',
       'yara_signatures_without_urls', 'yara_signatures_contentis_base64',
       'yara_signatures_Big_Numbers1', 'yara_signatures_with_urls',
       'yara_signatures_multiple_versions',
       'yara_signatures_invalid_trailer_structure',
       'yara_signatures_suspicious_packer_section', 'yara_signatures_vmdetect',
       'yara_signatures_Big_Numbers3', 'yara_signatures_Big_Numbers4',
       'yara_signatures_Big_Numbers2', 'yara_signatures_with_images',
       'yara_signatures_IsSuspicious', 'yara_signatures_invalid_xref_numbers',
       'yara_signatures_Big_Numbers0', 'yara_signatures_multiple_filtering',
       'yara_signatures_suspicious_launch_action',
       'yara_signatures_JBIG2_wrong_version', 'yara_signatures_header_evasion',
       'yara_signatures_MoleBoxv20',
       'yara_signatures_possible_includes_base64_packed_functions',
       'yara_signatures_FlateDecode_wrong_version',
       'yara_signatures_memory_shylock',
       'yara_signatures_maldoc_OLE_file_magic_number'],
      dtype='object')
```

Summary Stats

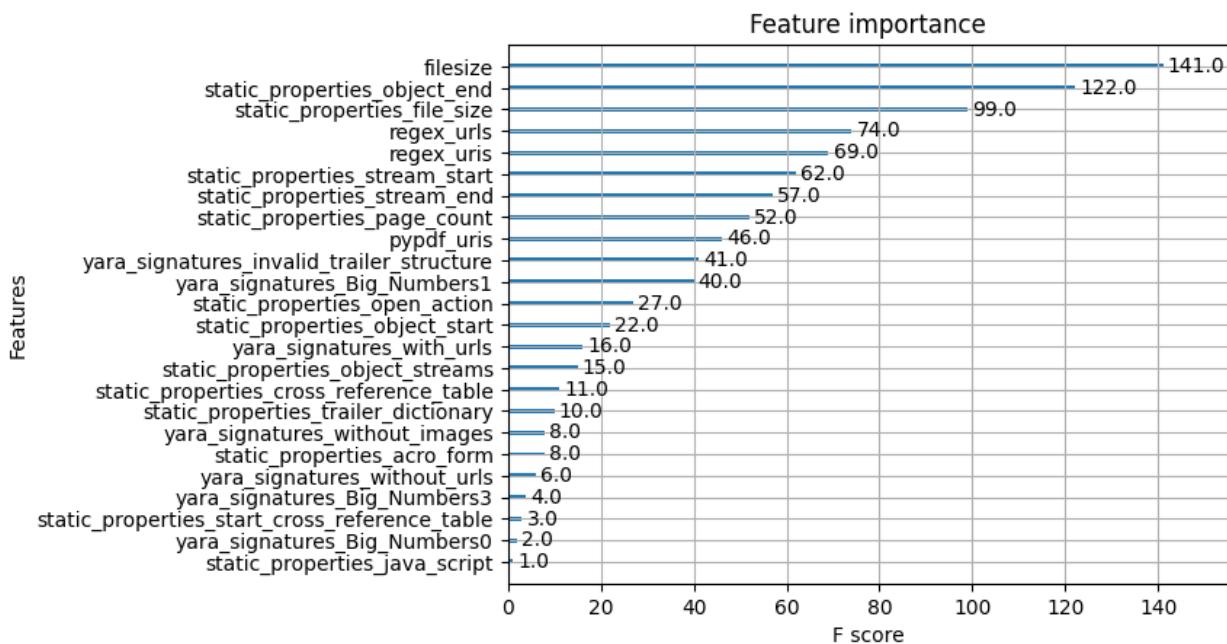
```
benign_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5000 entries, 0 to 4999
Data columns (total 54 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   filesize        5000 non-null    float64
 1   pypdf_uris     5000 non-null    int64  
 2   regex_uris     5000 non-null    int64  
 3   regex_urls     5000 non-null    int64  
 4   scripts_iframe 5000 non-null    int64  
 5   scripts_urls   5000 non-null    int64  
 6   static_properties_JBIG2Decode 4883 non-null    float64
 7   static_properties_XML_forms 4883 non-null    float64
 8   static_properties_acro_form 4883 non-null    float64
 9   static_properties_auto_action 4883 non-null    float64
 10  static_properties_colors 4883 non-null    float64
 11  static_properties_cross_reference_table 4883 non-null    float64
 12  static_properties_embedded_files 4883 non-null    float64
 13  static_properties_file_size 4984 non-null    float64
 14  static_properties_java_script 4883 non-null    float64
 15  static_properties_js 4883 non-null    float64
 16  static_properties_launch_action 4883 non-null    float64
 17  static_properties_object_end 4883 non-null    float64
 18  static_properties_object_start 4883 non-null    float64
 19  static_properties_object_streams 4883 non-null    float64
 20  static_properties_open_action 4883 non-null    float64
 21  static_properties_page_count 4883 non-null    float64
 22  static_properties_rich_media 4883 non-null    float64
 23  static_properties_start_cross_reference_table 4883 non-null    float64
 24  static_properties_stream_end 4883 non-null    float64
 25  static_properties_stream_start 4883 non-null    float64
 26  static_properties_trailer_dictionary 4883 non-null    float64
 27  yara_signatures_without_attachments 2420 non-null    object 
 28  yara_signatures_without_images 2381 non-null    object 
 29  yara_signatures_without_urls 809 non-null    object 
 30  yara_signatures_contentis_base64 2419 non-null    object 
 31  yara_signatures_Big_Numbers1 2208 non-null    object 
 32  yara_signatures_with_urls 1611 non-null    object 
 33  yara_signatures_multiple_versions 1052 non-null    object 
 34  yara_signatures_invalid_trailer_structure 851 non-null    object 
 35  yara_signatures_suspicious_packer_section 66 non-null    object 
 36  yara_signatures_vmdetect 199 non-null    object 
 37  yara_signatures_Big_Numbers3 45 non-null    object 
 38  yara_signatures_Big_Numbers4 6 non-null    object 
 39  yara_signatures_Big_Numbers2 13 non-null    object
```

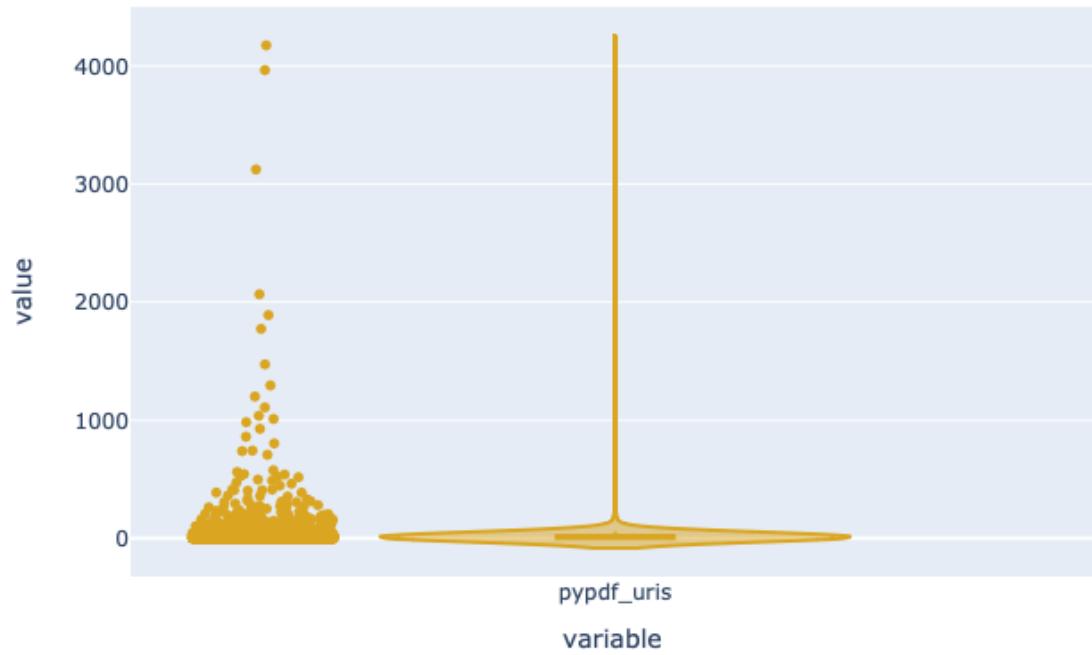
Similarly, data was also explored for the malicious dataframe. Further analysis was done by comparing them.

```
{'Verdict',
'yara_signatures_Big_Numbers4',
'yara_signatures_FlateDecode_wrong_version',
'yara_signatures_IsSuspicious',
'yara_signatures_JBIG2_wrong_version',
'yara_signatures_MoleBoxv20',
'yara_signatures_header_evasion',
'yara_signatures_maldoc_OLE_file_magic_number',
'yara_signatures_memory_shylock',
'yara_signatures_multiple_filtering'}
```

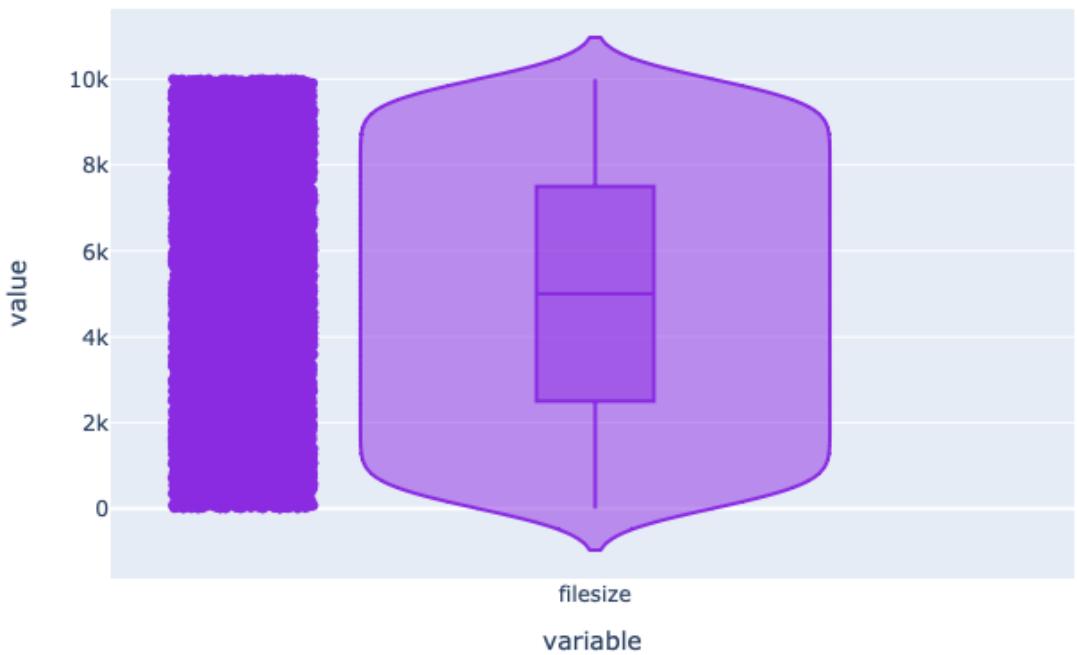
Correlation Analysis - Explore the correlation between features to identify potential relationships and dependencies within the dataset.



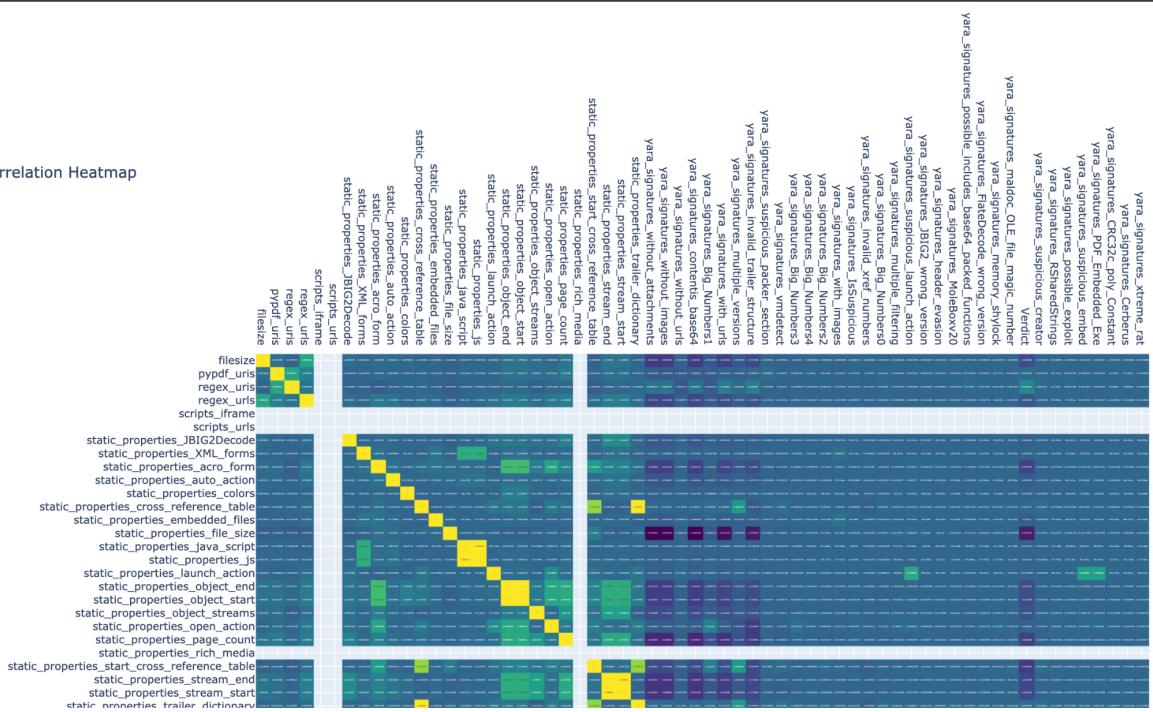
Violin plot for pypdf_uris Column

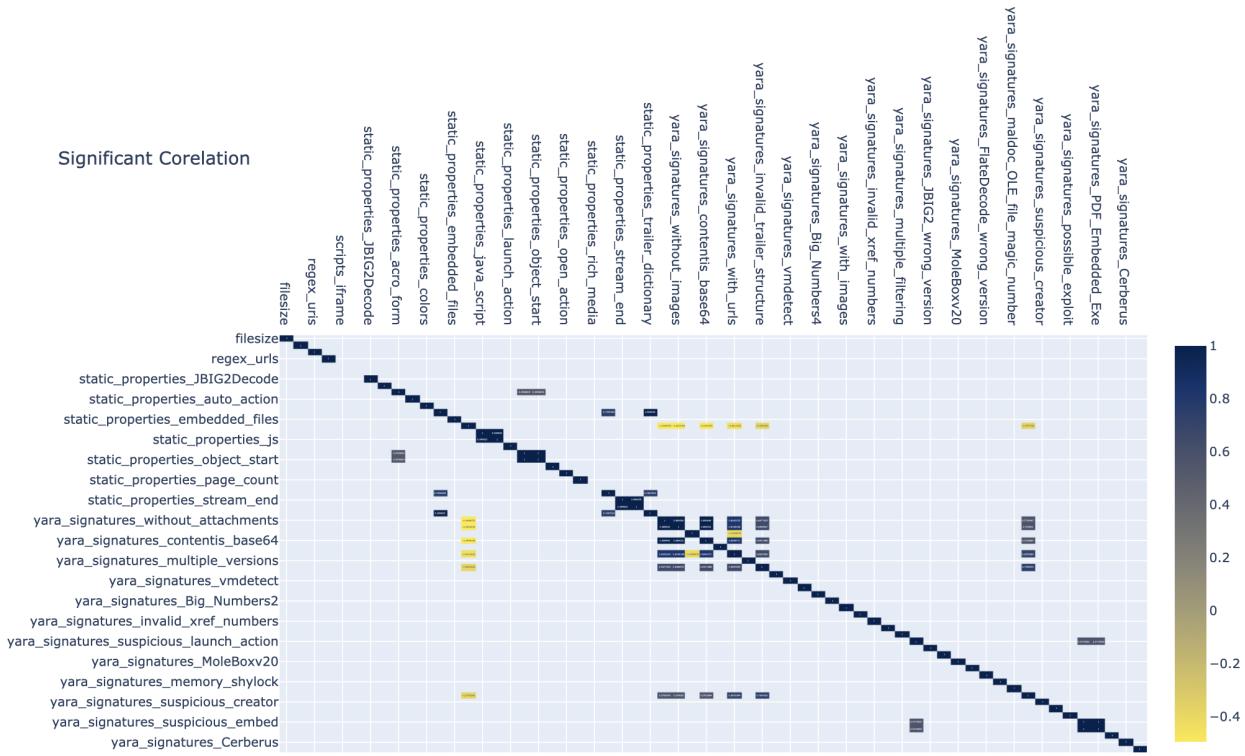


Violin plot for filesize Column



Correlation Heatmap





Visualization - Utilize visualizations to gain a deeper understanding of the dataset. This includes visualizing feature distributions, creating a correlation heatmap to highlight inter-feature relationships, examining class distributions, and identifying any emergent patterns.

Methodology

1. Using `predict_proba()`:

In the realm of classification models, relying solely on predictions (`classifier.predict(X)`) can be suboptimal. Instead, the use of `predict_proba(X)` is emphasized. This function provides the probability estimates for each class, offering a more nuanced understanding of the model's confidence in its predictions.

```
```python
Generate probability estimates for each class using predict_proba()
y_proba_base_xgb = xgb.predict_proba(X_test)
print(f'Shape of y_proba_base_xgb = {y_proba_base_xgb.shape} \n')
```

```

This code snippet showcases the utilization of `predict_proba()` with an XGBoost classifier (`xgb`). The resulting `y_proba_base_xgb` array contains probability estimates

for each class, allowing for a more comprehensive analysis of the model's predictive confidence.

ROC Curve, AUC, and Youden's J Statistic:

Metrics such as the Receiver Operating Characteristic (ROC) curve, Area Under the Curve (AUC), and Youden's J statistic play a pivotal role in evaluating the performance of classification models. These metrics provide insights into the trade-off between sensitivity and specificity and aid in selecting an optimal threshold for classification. The importance of these metrics is crucially acknowledged in model assessment.

1. ROC Curve and AUC:

The ROC curve is a fundamental tool for assessing binary classification models. It plots the true positive rate (sensitivity) against the false positive rate (1-specificity) across different classification thresholds. The Area Under the Curve (AUC) quantifies the overall performance of a model. A higher AUC indicates better discrimination between positive and negative instances. AUC is particularly useful in imbalanced datasets, providing a single, concise measure of a model's ability to distinguish between classes.

2. Youden's J Statistic:

Youden's J statistic, derived from the ROC curve, is a summary measure of a diagnostic test's performance. It is calculated using sensitivity and specificity values and is expressed as:

$$[J = \{Sensitivity\} + \{Specificity\} - 1]$$

This statistic aids in selecting an optimal classification threshold by maximizing the difference between sensitivity and specificity. In Python, Youden's J can be calculated using arrays for true positive rate (tpr) and false positive rate (fpr) obtained from an ROC curve:

```
```python
import numpy as np

Example values for true positive rate and false positive rate
tpr = np.array([0.2, 0.4, 0.6, 0.8, 1.0])
fpr = np.array([0.0, 0.1, 0.3, 0.5, 1.0])

Calculate Youden's J
youden_j = tpr + (1 - fpr) - 1
```

```

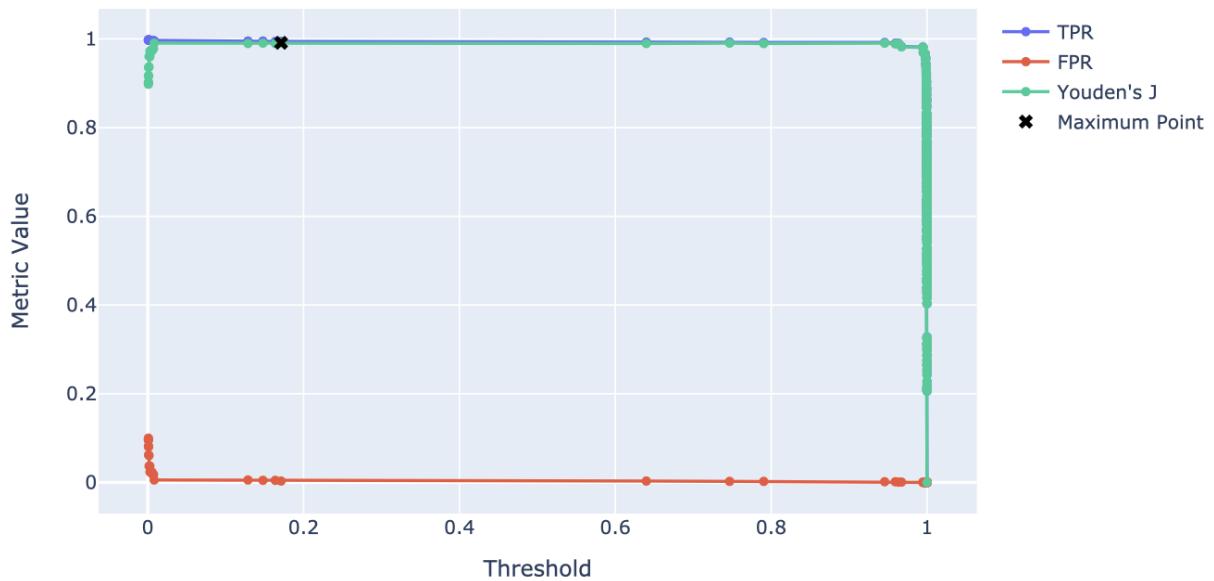
Find the index of the maximum J value
max_j_index = np.argmax(youden_j)

Print the maximum Youden's J value and corresponding coordinates
print("Maximum Youden's J:", youden_j[max_j_index])
print("Corresponding Sensitivity (True Positive Rate):", tpr[max_j_index])
print("Corresponding Specificity (1 - False Positive Rate):", 1 - fpr[max_j_index])
...

```



TPR and FPR for different thresholds



Understanding and utilizing these metrics are essential for making informed decisions about model performance, especially in scenarios where the consequences of false positives and false negatives are critical. They collectively provide a comprehensive evaluation of a model's discriminatory power and aid in selecting optimal thresholds for effective classification.

### Cross-Validation:

Cross-validation is a crucial step in model evaluation, playing a pivotal role in ensuring the robustness and generalizability of machine learning models. The process involves partitioning the dataset into multiple subsets, training the model on a subset, and evaluating its performance on the remaining subsets. This iterative approach aids in detecting issues related to overfitting and provides a more reliable estimate of a model's performance.

### Suggested Value of k in k-fold Cross-Validation:

The choice of the value of k, representing the number of folds in k-fold cross-validation, is essential for balancing validation efficiency and computational resources. A commonly suggested value is 10, which strikes a balance between providing a sufficiently large validation set and not excessively taxing computational resources.

### Python Implementation

In Python, the code snippet below demonstrates the application of k-fold cross-validation with an XGBoost classifier (`xgb`). The `cross\_val\_score` function evaluates the model using the specified number of folds, and the results are printed, including the mean accuracy across all folds.

```
```python
from sklearn.model_selection import cross_val_score

# X_train and y_train are the training data and labels
# skf_cv is the StratifiedKFold cross-validator

# Cross-Validation
cv_results = cross_val_score(xgb, X_train, y_train, cv=skf_cv, scoring='accuracy')
print("Cross-validation results:", cv_results)

# Calculate and print the mean accuracy across all folds
mean_accuracy = cv_results.mean()
print("Mean Accuracy:", mean_accuracy)
````
```

This Python code exemplifies the importance of cross-validation by evaluating an XGBoost classifier's performance using k-fold cross-validation and provides insights into the mean accuracy, offering a more comprehensive understanding of the model's generalization capabilities.

### LazyPredict: Rapid Model Selection and Prototyping

LazyPredict emerges as a powerful tool in the realm of machine learning, streamlining the model selection process and expediting the initial stages of model development. This automated approach offers quick insights into the performance of various classifiers, eliminating the need for extensive manual tuning.

#### Python Implementation:

In the following code snippet, LazyPredict is utilized to perform rapid model selection and prototype development. The dataset, represented by X and y, is split into training and testing sets. The LazyClassifier is employed with specific configurations, and the resulting models and predictions are printed for further analysis.

```

```python
# X_train, X_test, y_train, and y_test are the training and testing data and labels
# pdf_df is the DataFrame containing the dataset

# LazyPredict for Model Selection
X, y = pdf_df.drop('Verdict', axis=1), pdf_df['Verdict']

clf = LazyClassifier(verbose=0, ignore_warnings=True)
models, predictions = clf.fit(X_train, X_test, y_train, y_test)

# Print the list of models and their corresponding performance metrics
print(models)
```

```

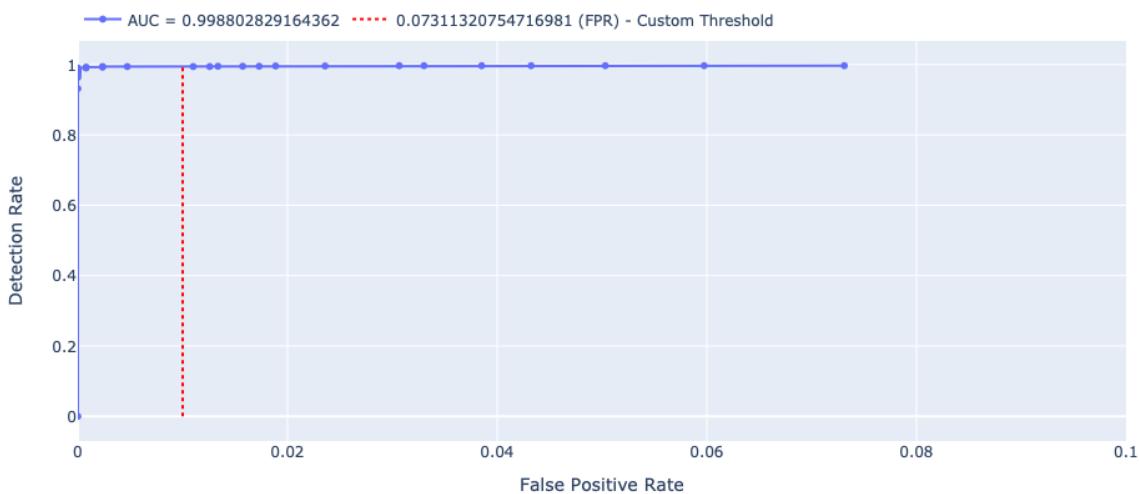
This Python code showcases the implementation of LazyPredict, providing a quick overview of various classifiers and their performance metrics. The tool's automation capabilities make it a valuable asset for efficiently exploring the landscape of potential models in the early stages of a machine learning project.

## Various Classifiers

### 1. ExtraTreesClassifier:

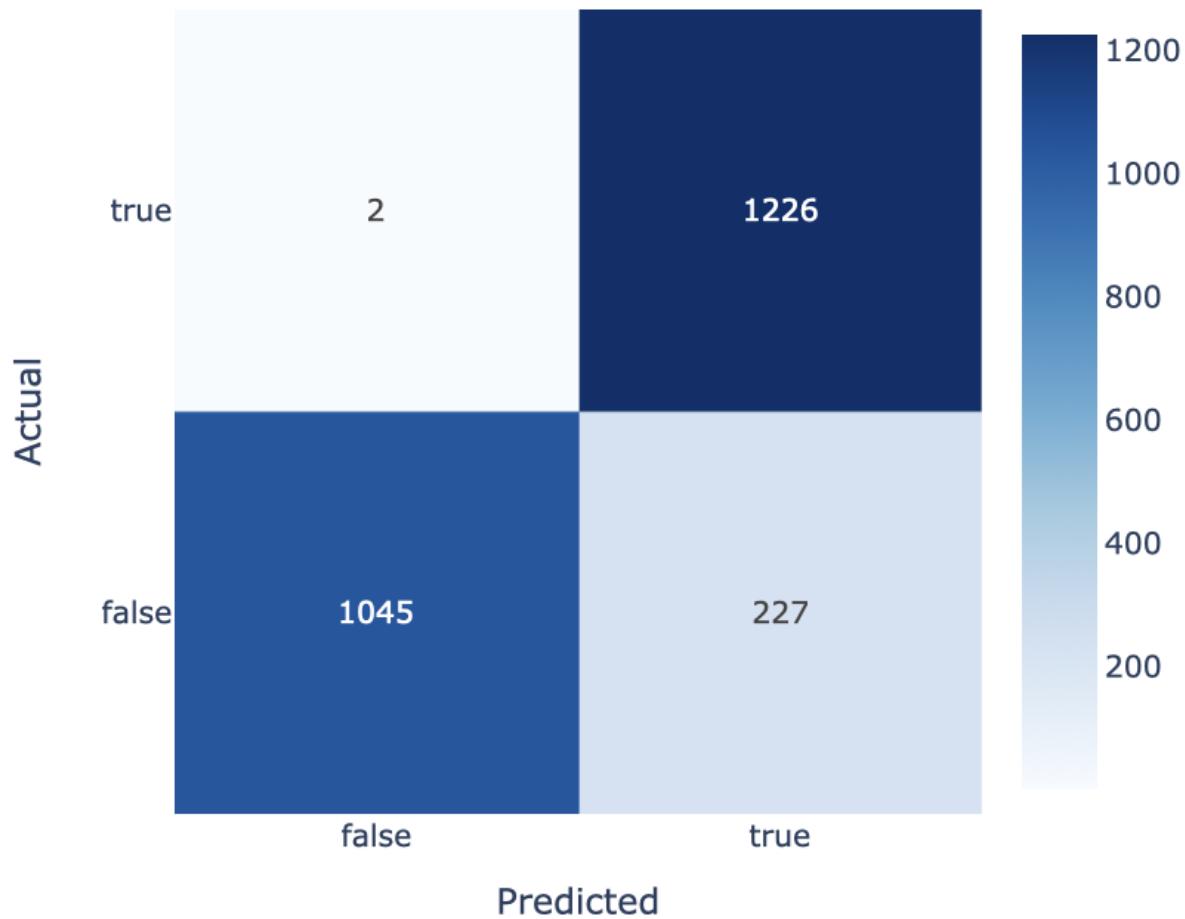
ExtraTreesClassifier is an ensemble learning method that builds a multitude of decision trees. Distinguishing itself by randomizing both feature and split selection, it enhances model robustness and is particularly effective for high-dimensional datasets and handling noisy data. Its strengths lie in the reduced risk of overfitting, high predictive accuracy, and robustness to outliers, though some parameter tuning may be required for optimal performance.

ROC Curve for Base ExtraTreesClassifier





### Confusion Matrix for Base etsc Model



---

#### 2. XGBClassifier (XGBoost):

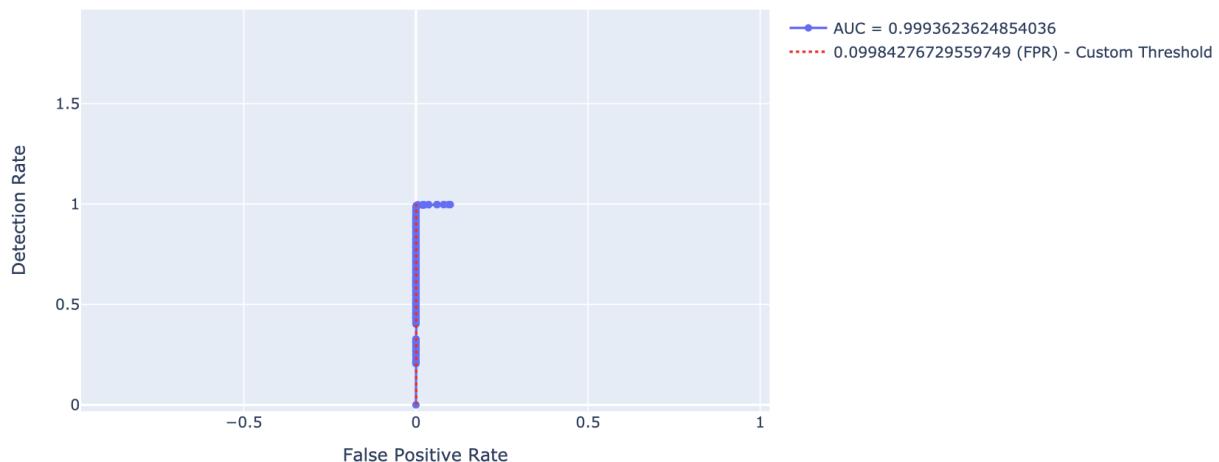
XGBClassifier, based on the XGBoost gradient boosting algorithm, stands out for its regularization techniques and efficient handling of missing data. It is renowned for its high predictive performance, effectiveness with diverse datasets, and the ability to analyze feature importance. While it delivers excellent results, careful parameter tuning is often necessary to harness its full potential.

```
▼ XGBClassifier
XGBClassifier(base_score=None, booster=None, callbacks=None,
 colsample_bylevel=None, colsample_bynode=None,
 colsample_bytree=None, device=None, early_stopping_rounds=None,
 enable_categorical=False, eval_metric=None, feature_types=None,
 gamma=None, grow_policy=None, importance_type=None,
 interaction_constraints=None, learning_rate=None, max_bin=None,
 max_cat_threshold=None, max_cat_to_onehot=None,
 max_delta_step=None, max_depth=None, max_leaves=None,
 min_child_weight=None, missing=nan, monotone_constraints=None,
 multi_strategy=None, n_estimators=None, n_jobs=None,
 num_parallel_tree=None, random_state=None, ...)
```

Shape of y\_proba = (2500, 2)

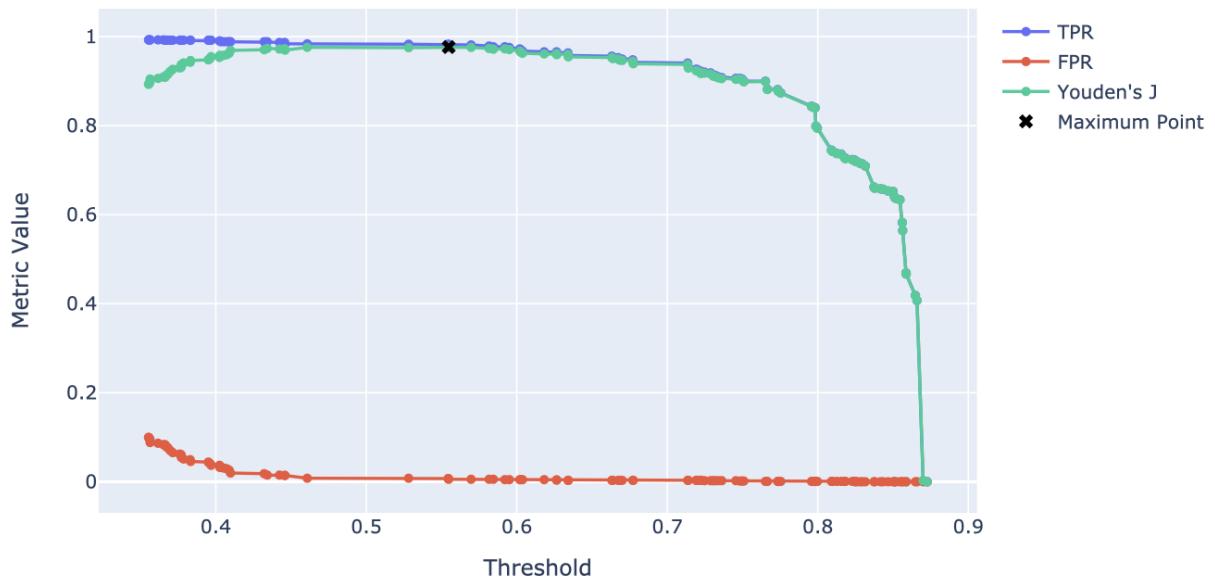
```
array([[1.33991241e-04, 9.99866009e-01],
 [9.99921679e-01, 7.83336873e-05],
 [9.99896407e-01, 1.03583116e-04],
 ...,
 [5.96046448e-06, 9.99994040e-01],
 [1.91926956e-05, 9.99980807e-01],
 [9.78936911e-01, 2.10631117e-02]], dtype=float32)
```

ROC Curve for Base XGB





### TPR and FPR for different Thresholds for XGBoost



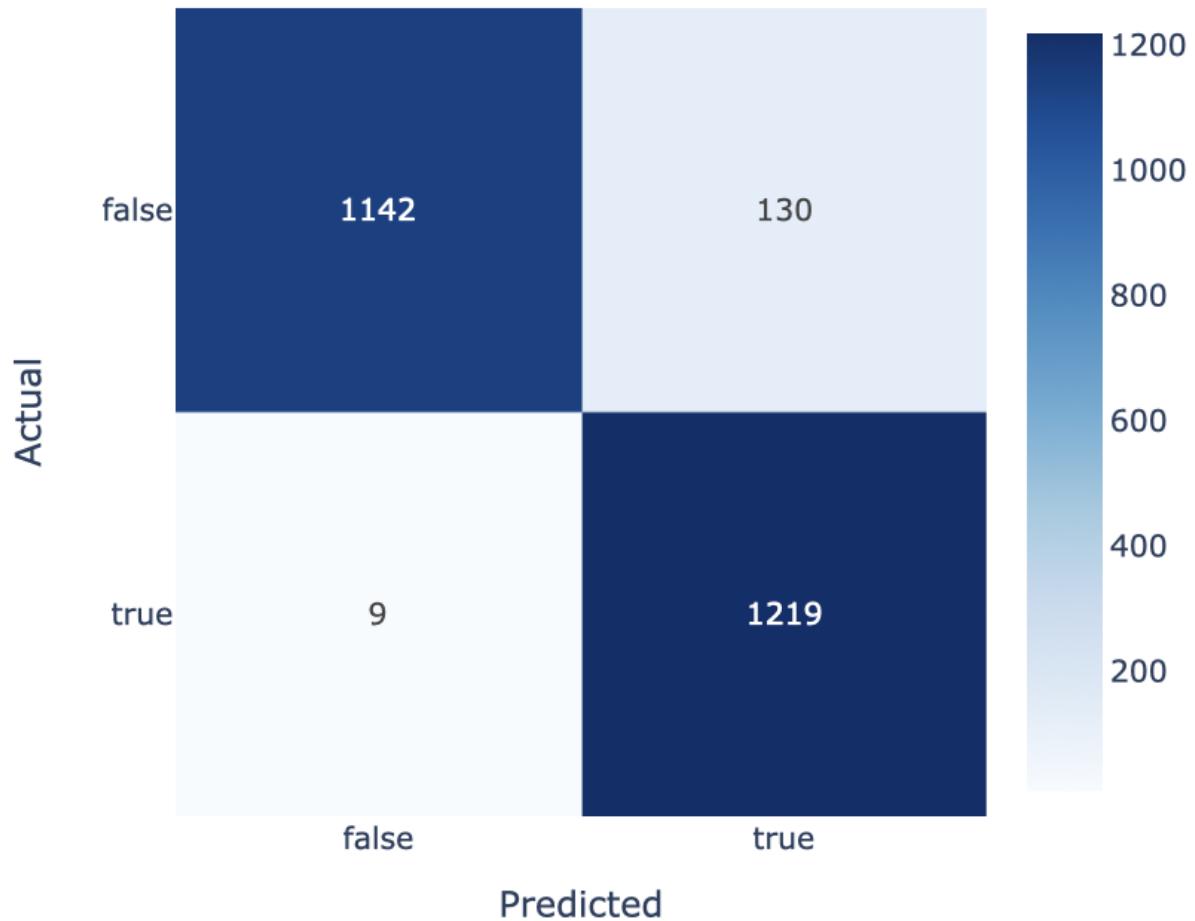
### 3. RandomForestClassifier:

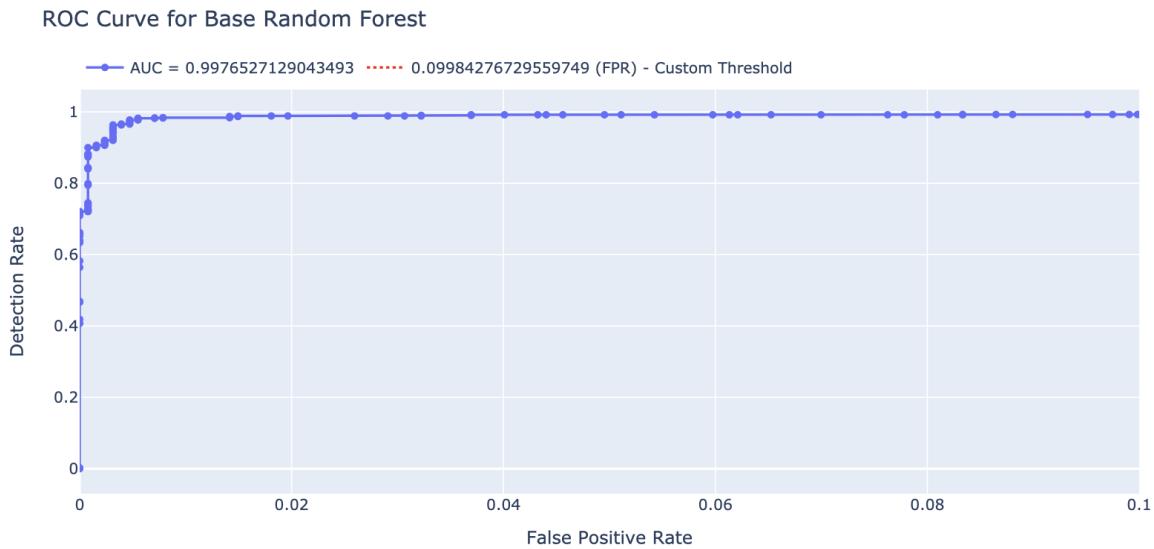
RandomForestClassifier is an ensemble method that leverages multiple decision trees. With random feature selection and bagging, it achieves robustness and is suitable for both classification and regression tasks. Resilient to overfitting, it excels in predictive accuracy and can handle large datasets with high dimensionality. However, its interpretability may be moderate compared to simpler models.





## Confusion Matrix for Base Random Forest Model





### **Extracting Features from JSON Files:**

#### **Downloading Features.zip from the Server:**

To initiate the feature extraction process, the features.zip file can be downloaded from the server. This task involves using the following commands:

```
```bash
# Downloading features.zip from the server
wget http://server-url/features.zip
```

```

**Additionally, to provide a visual guide, screenshots of the download process can be included.**

#### **Rationale Behind Choosing JSON as the Data Format:**

The choice of JSON as the data format for storing features is underpinned by its flexibility and human-readable structure. JSON (JavaScript Object Notation) offers a lightweight and easy-to-read format, making it well-suited for storing structured data. Its key-value pair representation allows for efficient organization of feature information. Moreover, JSON's compatibility with various programming languages ensures seamless integration into different environments. The human-readable nature of JSON facilitates easy inspection and debugging of feature data, contributing to the overall transparency of the feature extraction process.

#### **Using the Feature Extractor Tool:**

The feature extractor tool provided accepts three parameters: input path (`--in`), category (`--cat`), and output path (`--out`). For instance:

- Usage for malicious category:

```

```bash
extract_feature --in /dataset/malicious --cat malicious --out features
```
- Usage for benign category:
```bash
extract_feature --in /dataset/benign --cat benign --out features
```

```

The extracted features in JSON format are organized into respective categories within the output directory.

### **Already Extracted Features:**

Recognizing the computational intensity of feature extraction, a pre-extracted set of features is available on the server in the /features directory. Two options are provided for access:

1. Copy to Working Directory:

```

```bash
# Copying features from the server to the working directory
cp -r /features .
```

```

2. Download to Local Computer (Using scp):

```

```bash
# Downloading features using scp tool
scp -r cefil.joseph@plaksha.edu.in@10.1.46.56:/features .
```

```

This command facilitates the download of all features to the local computer using the Secure Copy Protocol (scp). The user needs to replace "cefil.joseph@plaksha.edu.in" with their specific credentials.

These features, available for direct use, expedite the feature extraction process and contribute to the efficiency of subsequent machine learning tasks.

### **Converting JSON Files into DataFrame and Memory Efficiency:**

Converting multiple JSON files into a DataFrame is a crucial step in preparing data for machine learning tasks. The following code illustrates a memory-efficient approach to loading and processing JSON files into a DataFrame:

```

```python
import os
import json
import pandas as pd

def load_dataset(dir_path):
    # Initialize an empty list to store individual DataFrames
    dfs = []
```

```

```

Iterate through each file in the folder
for filename in os.listdir(dir_path):
 if filename.endswith(".json"):
 file_path = os.path.join(dir_path, filename)

 # Read the JSON file into a dictionary
 with open(file_path, 'r') as file:
 json_data = json.load(file)
 reformatted_json = reformat_json(json_data, filename)

 # Convert the dictionary to a DataFrame
 df = pd.json_normalize(reformatted_json)

 # Append the DataFrame to the list
 dfs.append(df)

 # Concatenate all DataFrames in the list into a single DataFrame
 final_df = pd.concat(dfs, ignore_index=True)

return final_df
...

```

This function efficiently loads and processes JSON files by reading each file, converting its content into a DataFrame, and then concatenating these DataFrames into a single, coherent DataFrame. This approach is advantageous in terms of memory efficiency, especially when dealing with large datasets.

The use of `pd.json\_normalize()` ensures that nested JSON structures are appropriately flattened into the DataFrame, allowing for a clear and organized representation of the data. This memory-efficient conversion ensures optimal performance during subsequent stages of the machine learning pipeline.

## **Challenges Faced**

When extracting the benign files, we didn't encounter any issues. However, when using the same process to extract malicious files, we came across this error.

```
[29] benign_df = generate_dataframe(benign_json)
print(f'Shape = {benign_df.shape}\n')
benign_df.head()

Shape = (5000, 27)

 filesize pypdf_uris regex_uris regex_urls scripts_iframe scripts_urls static_properties_JBIG2Decode static_properties_XML_forms static_properties_acro_form static_properties
0 19303482.0 0 0 0 0 0 0 0 0
1 19303482.0 0 0 0 0 0 0 0 0
2 19303482.0 0 0 0 0 0 0 0 0
3 19303482.0 0 0 0 0 0 0 0 0
4 19303482.0 0 0 0 0 0 0 0 0

5 rows × 27 columns
```

```
[43] def get_json_file_list(directory_path):
 ...
 Function to retrieve the list of JSON file in a sub-directory.
 ...
 print('Parameters in get_json_file_list: ', locals())

 filename_list = []
 flag = 0
 i = 0

 for path in os.listdir(directory_path):
 # Check if current path is a file
 if not os.path.isfile(os.path.join(directory_path, path)):
 try:
 raise Exception('Unexpected behaviour - A directory is encountered!')
 except Exception as e:
 print(e)
 flag = 1
 return None
 else:
 #print(f'file {i}')
 i += 1
 filename_list.append(path)

 if flag == 0:
 return filename_list
```

```
def load_json_contents(dir_path):
 file_list = get_json_file_list(dir_path)
 print(f'Length of file_list = {len(file_list)}')
 json_list = []
 i = 0
 for filename in file_list:
 if i%1000 == 0:
 print(f'i = {i}, dir_path = {dir_path}')
 if filename.endswith('.json'):
 print(f'{filename} in file_list')
 i += 1
 filepath = os.path.join(dir_path, f)

 with open(filepath, 'r') as json_file:
 json_obj = json.load(json_file)
 reformatted_json = reformat_json(json_obj)
 json_list.append(reformatted_json)

 return json_list
```

```
[48] bt_list = get_json_file_list('./benign')
 mt_list = get_json_file_list('./malicious')
 len(set(bt_list).intersection(set(mt_list)))

Parameters in get_json_file_list: {'directory_path': './benign'}
Parameters in get_json_file_list: {'directory_path': './malicious'}
0
```

```
[33] os.listdir('./malicious')[0]
'3ecb3083e656ad062de81cf014fe5f2d5f6d4f2e0a3883e7ee05ae0b047000f1.properties.json'

[35] os.listdir('./malicious')[-1]
'2512f8935f3cd0ee5da8a786dc8bd0d8c78fef1a689c12f666d875a3b85f212f.properties.json'

[36] m_list = get_json_file_list('./malicious')
len(m_list)
5800
```

```
[52] malicious_json = load_json_contents('./malicious')
Parameters in get_json_file_list: {'directory_path': './malicious'}
Length of file_list = 5000
i = 0, dir_path = ./malicious
True

FileNotFoundError: [Errno 2] No such file or directory: './malicious/44a5ad451831a30c52b50c4f27feffedfc5de0b0630c48512b1c71bbfffc5d9f0.properties.json'
----- in cell line: 1{}()
----> 1 malicious_json = load_json_contents('./malicious')

[53] malicious_json = load_json_contents(dir_path)
 12 filepath = os.path.join(dir_path, f)
 13
> 14 with open(filepath, 'r') as json_file:
 15 json_obj = json.load(json_file)
 16 reformatted_json = reformat_json(json_obj)

FileNotFoundError: [Errno 2] No such file or directory: './malicious/44a5ad451831a30c52b50c4f27feffedfc5de0b0630c48512b1c71bbfffc5d9f0.properties.json'

SEARCH STACK OVERLOW
```

To overcome this we tried converting from int to float

```
def reformat_json(json_dict):
 reformatted_json = {
 'filesize': float(json_dict['file_size']),
 'pypdf_uris': len(json_dict['pypdf_uris']),
 'regex_uris': len(json_dict['regex_uris']),
 'regex_uris': len(json_dict['regex_uris']),
 'scripts_iframe': len(json_dict['scripts']['iframe']),
 'scripts_urls': len(json_dict['scripts']['urls']),
 }

 # Process static_properties
 for key, value in json_dict['static_properties'].items():
 reformatted_json[f'static_properties_{key}'] = int(value)

 # Process yara_signatures
 for signature in json_dict['yara_signatures']:
 reformatted_json[f'yara_signatures_{signature}'] = True

 return reformatted_json
```

```

Initialize an empty list to store individual DataFrames
dfs = []

Iterate through each file in the folder
for filename in os.listdir(folder_path):
 if filename.endswith(".json"):
 file_path = os.path.join(folder_path, filename)

 # Read the JSON file into a dictionary
 with open(file_path, 'r') as file:
 json_data = json.load(file)
 reformatted_json = reformat_json(json_data)

 # Convert the dictionary to a DataFrame
 df = pd.json_normalize(reformatted_json)

 # Append the DataFrame to the list
 dfs.append(df)

Concatenate all DataFrames in the list into a single DataFrame
final_df = pd.concat(dfs, ignore_index=True)

Print or use the final DataFrame as needed
final_df.head()

```

ValueError: invalid literal for int() with base 10: "3.6"

```

Process static_properties
for key, value in json_dict['static_properties'].items():
 reformatted_json[f'static_properties_{key}'] = float(value)

```

```

Initialize an empty list to store individual DataFrames
dfs = []

Iterate through each file in the folder
for filename in os.listdir(folder_path):
 if filename.endswith(".json"):
 file_path = os.path.join(folder_path, filename)

 # Read the JSON file into a dictionary
 with open(file_path, 'r') as file:
 json_data = json.load(file)
 reformatted_json = reformat_json(json_data)

 # Convert the dictionary to a DataFrame
 df = pd.json_normalize(reformatted_json)

 # Append the DataFrame to the list
 dfs.append(df)

Concatenate all DataFrames in the list into a single DataFrame
final_df = pd.concat(dfs, ignore_index=True)

Print or use the final DataFrame as needed
final_df.head()

```

|   | filesize | pydf_uris | regex_uris | regex_uris | scripts_iframe | scripts_uris | static_properties_MIMEDecode | static_properties_XMLForms | static_properties_acro_form | static_properties_auto_action | ... | yara_signatures_Big_Numbers |
|---|----------|-----------|------------|------------|----------------|--------------|------------------------------|----------------------------|-----------------------------|-------------------------------|-----|-----------------------------|
| 0 | 375641.0 | 3         | 0          | 5          | 0              | 0            | 0.0                          | 0.0                        | 0.0                         | 0.0                           | ... | NaN                         |
| 1 | 187630.0 | 0         | 0          | 5          | 0              | 0            | 0.0                          | 0.0                        | 0.0                         | 0.0                           | ... | NaN                         |
| 2 | 365605.0 | 0         | 0          | 2          | 0              | 0            | 0.0                          | 0.0                        | 0.0                         | 0.0                           | ... | NaN                         |
| 3 | 331367.0 | 32        | 0          | 8          | 0              | 0            | 0.0                          | 0.0                        | 0.0                         | 0.0                           | ... | NaN                         |
| 4 | 141494.0 | 3         | 0          | 2          | 0              | 0            | 0.0                          | 0.0                        | 0.0                         | 0.0                           | ... | NaN                         |

5 rows × 53 columns

```

malicious = load_dataset("./malicious", "Malicious")
malicious.head()

```

ValueError: invalid literal for int() with base 10: '612fdd72a7ca'

```

Process static properties
for key, value in json_dict['static_properties'].items():
 reformatted_json[f'static_properties_{key}'] = float(value)

Process yara_signatures

```

ValueError: could not convert string to float: '-'

SEARCH STACK OVERFLOW

We found that there was a hyphen somewhere which interrupted the entire process.

```
✓ [113] def reformat_json(json_dict, filename):
 reformatted_json = {'filesize': float(json_dict['file_size']),
 'pypdf_uris': len(json_dict['pypdf_uris']),
 'regex_uris': len(json_dict['regex_uris']),
 'regex_urls': len(json_dict['regex_urls']),
 'scripts_iframe': len(json_dict['scripts']['iframe']),
 'scripts_urls': len(json_dict['scripts']['urls']),
 }

 # Process static_properties
 for key, value in json_dict['static_properties'].items():
 try:
 reformatted_json[f'static_properties_{key}'] = float(value)
 except ValueError as ve:
 print(ve)
 print('key = ', key)
 print(filename, '\n')

 # Process yara_signatures
 for signature in json_dict['yara_signatures']:
 reformatted_json[f'yara_signatures_{signature}'] = True

 return reformatted_json
```

```
malicious_df = load_dataset('../malicious', 'Malicious')
malicious_df.head()
could not convert string to float: '-'
key = file_size
uc233897d712ffac1eda515101ebec61644a088359c1f1m988c539f8e411el.properties.json

 filesize pypdf_uris regex_uris regex_urls scripts_iframe scripts_uris static_properties_38IG2Decade static_properties_386_forms static_properties_acro_form static_properties_auto_action ... yara_signatures_suspicious_law
0 46292.0 22 22 42 0 0 0.0 0.0 0.0 0.0 ...
1 36077.0 22 22 22 0 0 0.0 0.0 0.0 0.0 ...
2 246234.0 21 28 27 0 0 0.0 0.0 0.0 0.0 ...
3 148759.0 24 24 43 0 0 0.0 0.0 0.0 0.0 ...
4 35035.0 22 22 22 0 0 0.0 0.0 0.0 0.0 ...

5 rows × 52 columns
```

Now that we found the issue, we modified the code to tackle the same. We tried finding the file to see which file had the issue. Once we found the root cause, we made the necessary modifications to the code. Hence, we were able to circumvent this issue.

```

[✓] [0s] import numpy as np

[✓] [0s] def reformat_json(json_dict, filename):
 reformatted_json = {'filesize': float(json_dict['file_size']),
 'pypdf_uris': len(json_dict['pypdf_uris']),
 'regex_uris': len(json_dict['regex_uris']),
 'regex_urls': len(json_dict['regex_urls']),
 'scripts_iframe': len(json_dict['scripts']['iframe']),
 'scripts_urls': len(json_dict['scripts']['urls']),
 }

 # Process static_properties
 for key, value in json_dict['static_properties'].items():
 try:
 reformatted_json[f'static_properties_{key}'] = float(value)
 except ValueError as ve:
 print(ve)
 print('key = ', key)
 print(filename, '\n')
 reformatted_json[f'static_properties_{key}'] = np.NaN

 # Process yara_signatures
 for signature in json_dict['yara_signatures']:
 reformatted_json[f'yara_signatures_{signature}'] = True

 return reformatted_json

```

```

[✓] [128] malicious_df = load_dataset('../malicious', 'Malicious')
[✓] [128] malicious_df.head()
could not convert string to float: '-'
key = filesize=46292.0,filename=0c333897d712fac1e1da515101ebef01644a088359c1f1ad98c539f4e411el.properties.json
 filesize pypdf_uris regex_uris regex_urls scripts_iframe scripts_uris static_properties_0BIG20encode static_properties_20LForms static_properties_acroForm static_properties_autoAction ... yara_signatures_suspicious_law
0 46292.0 22 22 42 0 0 0.0 0.0 0.0 0.0 ...
1 36077.0 22 22 22 0 0 0.0 0.0 0.0 0.0 ...
2 246234.0 21 21 27 0 0 0.0 0.0 0.0 0.0 ...
3 146759.0 24 24 43 0 0 0.0 0.0 0.0 0.0 ...
4 35035.0 22 22 22 0 0 0.0 0.0 0.0 0.0 ...
5 rows × 52 columns
[✓] [128] malicious_df['static_properties_file_size'].isna().sum()
3

```

This was the primary challenge we came across. Moreover, it is to be noted that there were more obvious challenges we faced in building and training the models. For instance, during hyperparameter tuning and maintaining the balance between false positive rates and detection rates.

## Model Evaluation

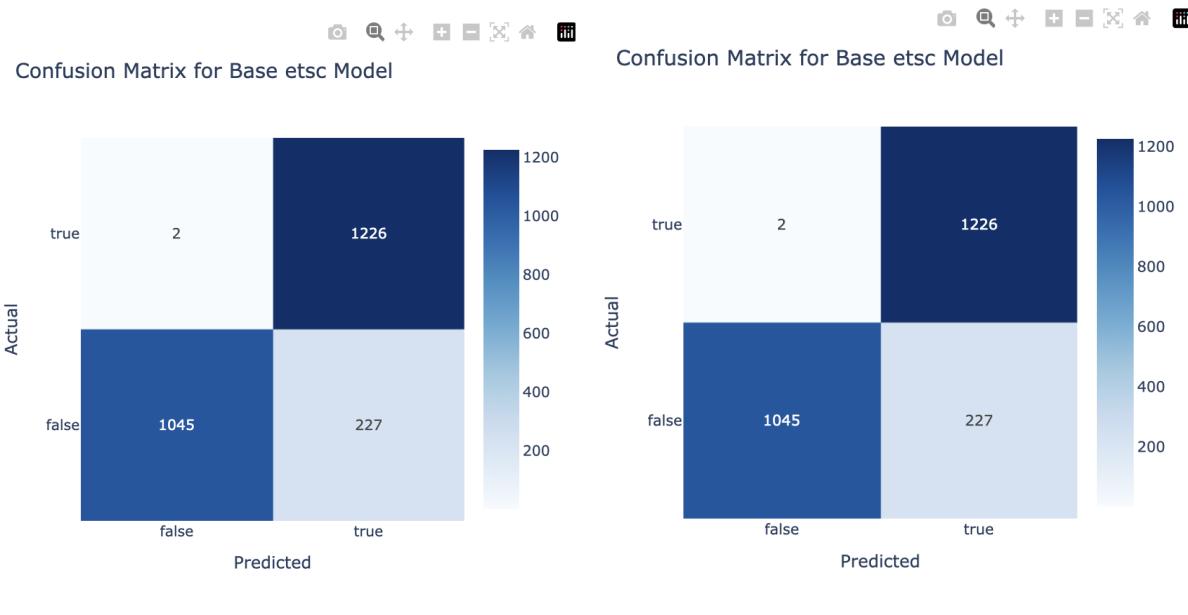
In the subsequent sections of this report, detailed results and insights from the model evaluation phase will be presented, providing a comprehensive overview of the model's performance and its optimization journey through hyperparameter tuning and feature importance analysis.

| Model                         |          |          |          |          |  |
|-------------------------------|----------|----------|----------|----------|--|
| LGBMClassifier                | 0.995600 | 0.995578 | 0.995578 | 0.995600 |  |
| XGBClassifier                 | 0.995600 | 0.995578 | 0.995578 | 0.995600 |  |
| RandomForestClassifier        | 0.995600 | 0.995578 | 0.995578 | 0.995600 |  |
| ExtraTreesClassifier          | 0.995600 | 0.995563 | 0.995563 | 0.995600 |  |
| BaggingClassifier             | 0.994000 | 0.993963 | 0.993963 | 0.994000 |  |
| AdaBoostClassifier            | 0.992800 | 0.992770 | 0.992770 | 0.992800 |  |
| DecisionTreeClassifier        | 0.992400 | 0.992377 | 0.992377 | 0.992400 |  |
| LabelPropagation              | 0.990400 | 0.990383 | 0.990383 | 0.990400 |  |
| LabelSpreading                | 0.990400 | 0.990383 | 0.990383 | 0.990400 |  |
| ExtraTreeClassifier           | 0.989600 | 0.989625 | 0.989625 | 0.989600 |  |
| KNeighborsClassifier          | 0.988800 | 0.988782 | 0.988782 | 0.988800 |  |
| PassiveAggressiveClassifier   | 0.983600 | 0.983672 | 0.983672 | 0.983601 |  |
| SVC                           | 0.983200 | 0.983110 | 0.983110 | 0.983198 |  |
| LinearSVC                     | 0.975600 | 0.975740 | 0.975740 | 0.975602 |  |
| LogisticRegression            | 0.975600 | 0.975698 | 0.975698 | 0.975602 |  |
| CalibratedClassifierCV        | 0.972800 | 0.972989 | 0.972989 | 0.972802 |  |
| SGDClassifier                 | 0.972800 | 0.972848 | 0.972848 | 0.972801 |  |
| RidgeClassifierCV             | 0.965600 | 0.965815 | 0.965815 | 0.965602 |  |
| LinearDiscriminantAnalysis    | 0.965600 | 0.965801 | 0.965801 | 0.965602 |  |
| RidgeClassifier               | 0.965200 | 0.965408 | 0.965408 | 0.965202 |  |
| GaussianNB                    | 0.949600 | 0.950063 | 0.950063 | 0.949587 |  |
| QuadraticDiscriminantAnalysis | 0.948800 | 0.949376 | 0.949376 | 0.948772 |  |
| BernoulliNB                   | 0.943600 | 0.944294 | 0.944294 | 0.943548 |  |
| Perceptron                    | 0.932800 | 0.932568 | 0.932568 | 0.932776 |  |
| NuSVC                         | 0.922800 | 0.923065 | 0.923065 | 0.922802 |  |
| NearestCentroid               | 0.828800 | 0.831522 | 0.831522 | 0.824947 |  |
| DummyClassifier               | 0.491200 | 0.500000 | 0.500000 | 0.323602 |  |

| Model                         | Time Taken |
|-------------------------------|------------|
| LGBMClassifier                | 0.339455   |
| XGBClassifier                 | 0.302155   |
| RandomForestClassifier        | 0.793701   |
| ExtraTreesClassifier          | 0.660622   |
| BaggingClassifier             | 0.398538   |
| AdaBoostClassifier            | 0.757857   |
| DecisionTreeClassifier        | 0.110652   |
| LabelPropagation              | 5.041122   |
| LabelSpreading                | 8.359449   |
| ExtraTreeClassifier           | 0.045881   |
| KNeighborsClassifier          | 0.334909   |
| PassiveAggressiveClassifier   | 0.086584   |
| SVC                           | 0.741627   |
| LinearSVC                     | 1.221125   |
| LogisticRegression            | 0.727985   |
| CalibratedClassifierCV        | 3.258137   |
| SGDClassifier                 | 0.168839   |
| RidgeClassifierCV             | 0.115328   |
| LinearDiscriminantAnalysis    | 0.283340   |
| RidgeClassifier               | 0.057635   |
| GaussianNB                    | 0.051902   |
| QuadraticDiscriminantAnalysis | 0.106978   |
| BernoulliNB                   | 0.067346   |
| Perceptron                    | 0.071212   |
| NuSVC                         | 7.435537   |
| NearestCentroid               | 0.118990   |
| DummyClassifier               | 0.047612   |

## Confusion Matrix

Analyzing the model's performance using a confusion matrix provides a detailed breakdown of its predictions. This matrix helps in understanding true positives, true negatives, false positives, and false negatives, providing insights into the model's strengths and weaknesses.



## Accuracy Score

The accuracy score is a fundamental metric that evaluates the overall correctness of the model's predictions. It is calculated as the ratio of correct predictions to the total number of predictions made by the model.

## ROC Curve

Visualizing the model's performance across different probability thresholds is crucial for understanding its sensitivity to false positive rates. In this project, a meticulous approach is taken to choose threshold values in models that produce a false positive (FP) rate of 0.1% or lower. Among these, the model with the highest detection rate (DR) is selected. The chosen threshold, along with the corresponding DR and FP, is then reported. This information is pivotal in determining the best project award winners.

```

Print the Detection Rate of Malicious Files along with corresponding
False Positive Rate and Threshold values
Loading...
print(f'DR vs FP with Thresholds')
for detection_rate, false_positive_rate, thr in zip(tpr, fpr, thresholds):
 print(f'{100 * detection_rate:>6.3f}%', f'{10 * false_positive_rate:10.4f}%', f'{thr:15.4f}')

DR vs FP with Thresholds
0.000% 0.0000% 1.0000
20.603% 0.0000% 1.0000
20.684% 0.0000% 1.0000
20.928% 0.0000% 1.0000
21.091% 0.0000% 1.0000
21.254% 0.0000% 1.0000
21.498% 0.0000% 1.0000
21.743% 0.0000% 1.0000
22.720% 0.0000% 1.0000
24.267% 0.0000% 1.0000
25.000% 0.0000% 1.0000
25.489% 0.0000% 1.0000
26.303% 0.0000% 1.0000
26.547% 0.0000% 1.0000
26.710% 0.0000% 1.0000
27.362% 0.0000% 1.0000
27.524% 0.0000% 1.0000
28.583% 0.0000% 1.0000
28.746% 0.0000% 1.0000
29.805% 0.0000% 1.0000
30.049% 0.0000% 1.0000
30.456% 0.0000% 1.0000
30.945% 0.0000% 1.0000
31.189% 0.0000% 1.0000
31.270% 0.0000% 1.0000
31.433% 0.0000% 1.0000
32.410% 0.0000% 1.0000
32.573% 0.0000% 1.0000
32.818% 0.0000% 1.0000
40.309% 0.0000% 1.0000
40.391% 0.0000% 1.0000
41.612% 0.0000% 1.0000
41.694% 0.0000% 1.0000
42.508% 0.0000% 1.0000
42.834% 0.0000% 1.0000
43.078% 0.0000% 1.0000

```

## Hyperparameter Tuning

Experimenting with hyperparameter tuning is an essential step in optimizing model performance. By systematically adjusting hyperparameters and assessing their impact on the model's accuracy, the goal is to find the most effective configuration for achieving optimal results.

## Feature Importance

Investigating feature importance is crucial for understanding which features contribute most significantly to the model's predictions. This analysis aids in refining the model, identifying key factors that influence its decision-making process.

## Ethical Considerations

- 1. Informed Consent:**  
Ensure that individuals contributing to the project, especially those involved in providing datasets or data samples, are fully informed about the nature of the project and the potential risks associated with handling malicious content.
- 2. Privacy and Anonymity:**  
Prioritize user privacy by anonymizing any sensitive information within the dataset. Avoid collecting or storing personally identifiable information (PII) unless absolutely necessary for the project's objectives.
- 3. Data Security:**  
Implement robust security measures to safeguard against potential data breaches. Malicious content should be handled securely to prevent unintended exposure or misuse.
- 4. Legal Compliance:**  
Adhere to relevant laws and regulations governing the handling of sensitive and potentially harmful content. This includes compliance with data protection regulations and laws related to cybersecurity.
- 5. Secure Development Practices:**  
Employ secure coding practices to minimize the risk of vulnerabilities that could be exploited by malicious content. Regularly update and patch software to ensure the latest security measures are in place.
- 6. Ethical Use of Findings:**  
Clearly define the purpose of the project and ensure that any findings or results obtained from the analysis of malicious content are used ethically. Avoid any use that could harm individuals, organizations, or society.
- 7. Transparency:**  
Maintain transparency in the project's goals, methodologies, and potential impact. Provide clear explanations of how malicious content is being handled and the precautions taken to mitigate risks.
- 8. Collaboration with Authorities:**  
In cases where the project involves handling potentially illegal content, collaborate with relevant authorities and law enforcement agencies. Report any findings that may indicate criminal activity.
- 9. Bias and Fairness:**  
Be aware of and address biases that may arise in the handling and analysis of malicious content. Ensure fairness in the evaluation and reporting of results.

**10. End-User Safety:**

Prioritize the safety of end-users who may be indirectly affected by the project.  
Implement measures to prevent unintended consequences or harm resulting from the handling of malicious content.

**11. Continuous Ethical Review:**

Regularly review and reassess the ethical implications of the project as it progresses.  
Adapt ethical considerations based on new insights, changing circumstances, or emerging ethical standards.

On addressing these ethical considerations, this project can contribute to cybersecurity advancements while upholding the highest standards of ethical conduct and responsibility.